

Object Detection (Plus some bonuses)

EECS 442 – Jeong Joon Park
Winter 2024, University of Michigan

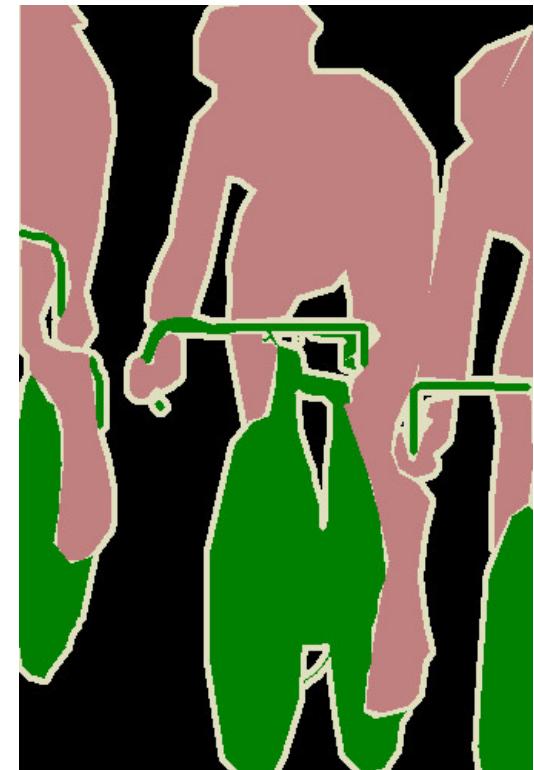
Last Time

“Semantic Segmentation”: Label each pixel with the object category it belongs to.

Input



Target



Today – Object Detection

“Object Detection”: Draw a box around each instance of a list of categories

Input

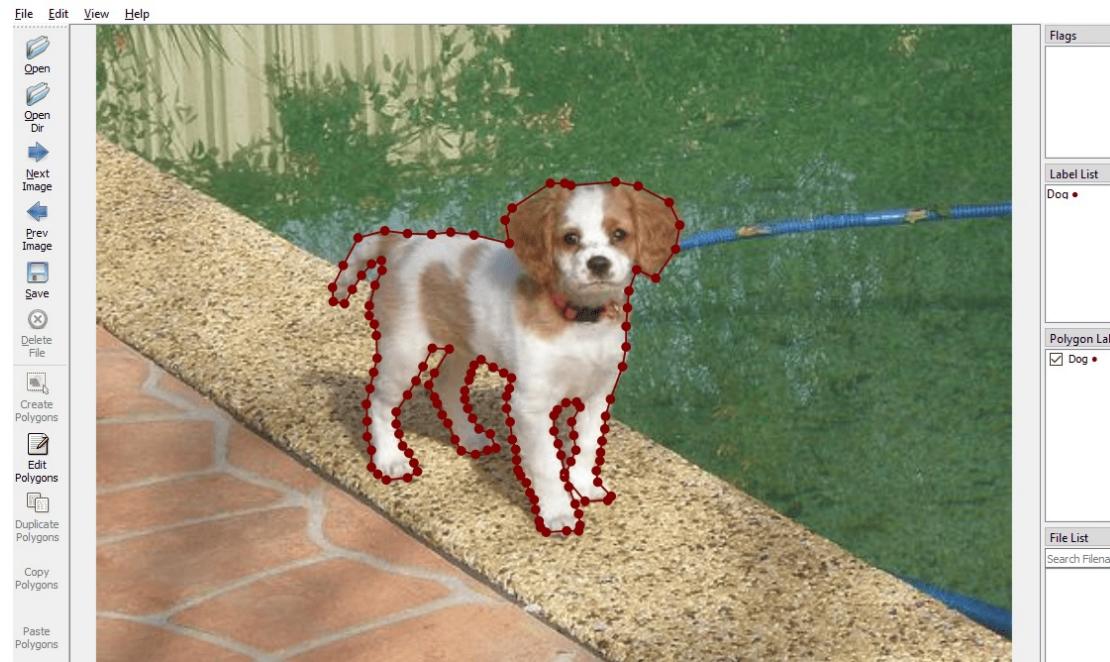


Target

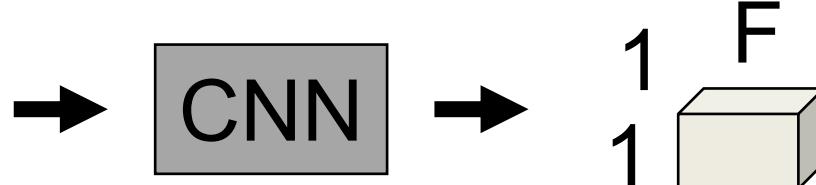


Why Object Detection?

- Instance segmentation?
- Efficiency: real-time applications
- Data annotation requirement



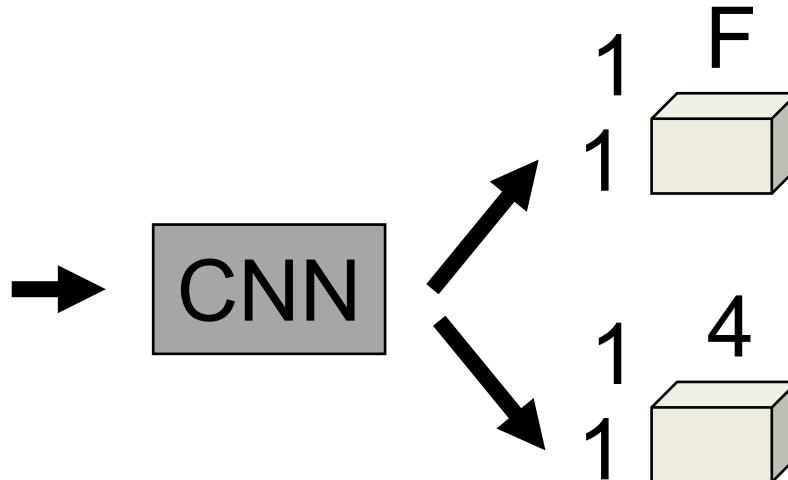
The Wrong Way To Do It



Starting point:

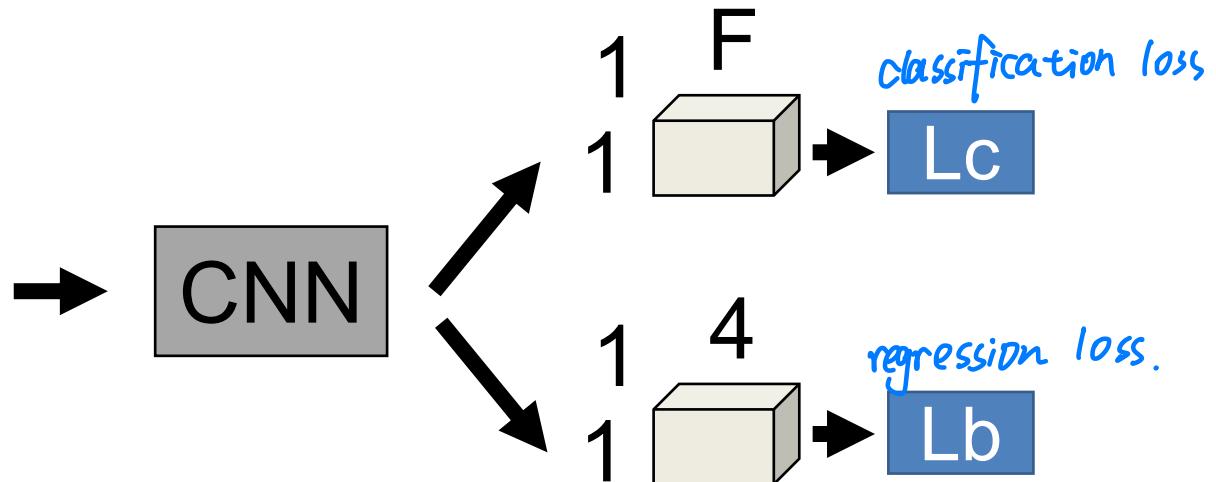
Can predict the probability of F classes
 $P(\text{cat}), P(\text{goose}), \dots P(\text{tractor})$

The Wrong Way To Do It



Add another output (why not):
Predict the *bounding box* of the object
[x,y,width,height] or [minX,minY,maxX,maxY]

The Wrong Way To Do It



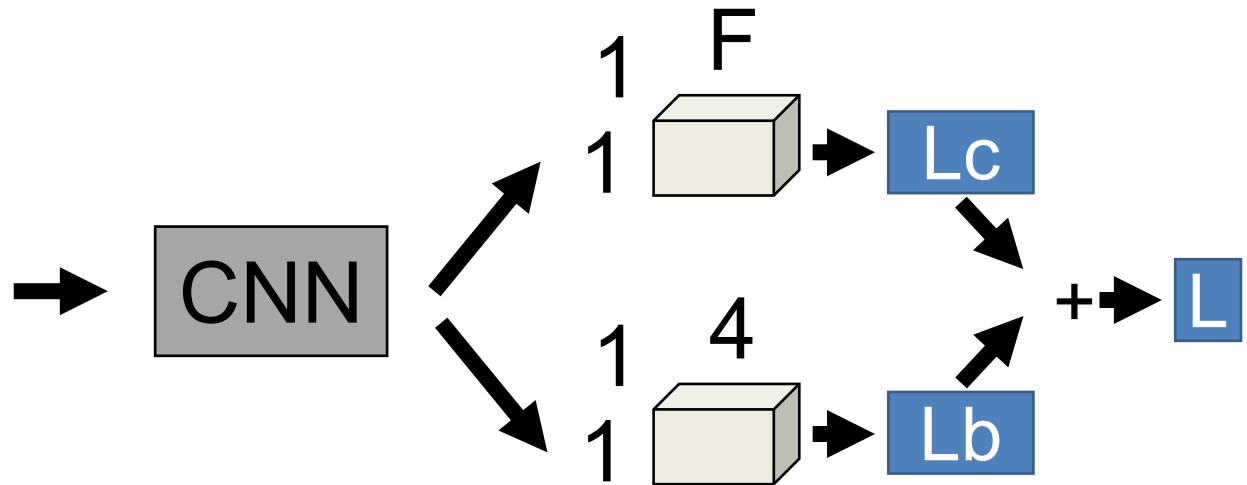
Put a loss on it:

Penalize mistakes on the classes with

L_c = negative log-likelihood

L_b = L2 loss

The Wrong Way To Do It

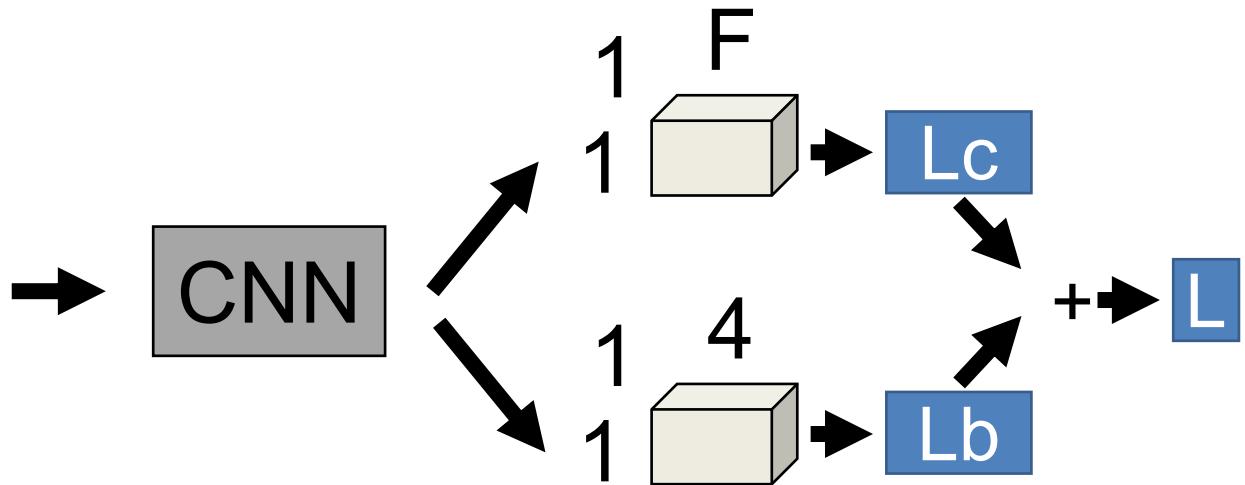


Add losses, backpropagate

$$\text{Final loss: } L = L_c + \lambda L_b$$

Why do we need the λ ?

The Wrong Way To Do It

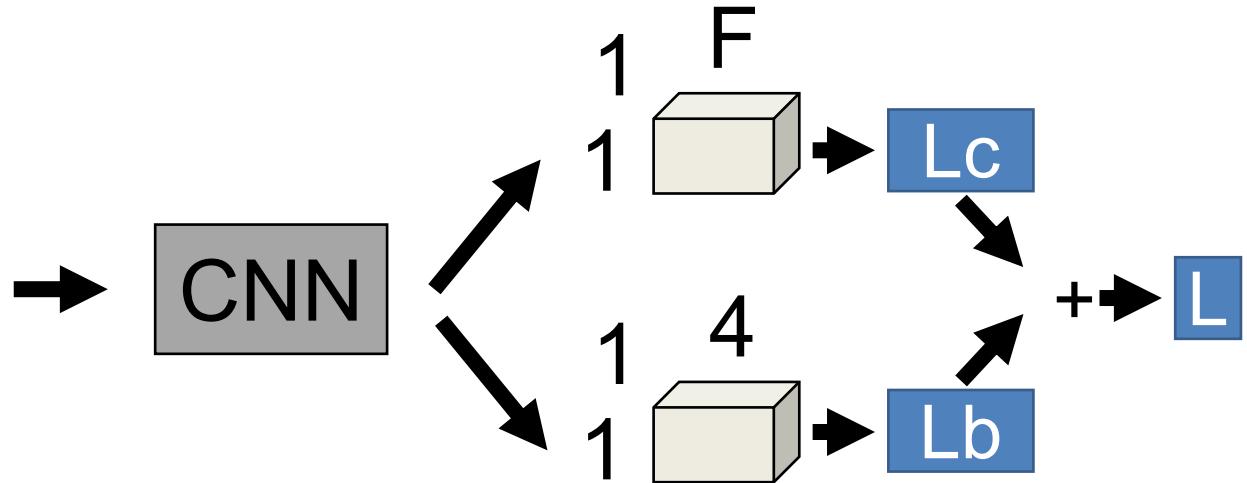


Now there are two ducks.

How many outputs do we need?

$$F, 4, F, 4 = 2*(F+4)$$

The Wrong Way To Do It

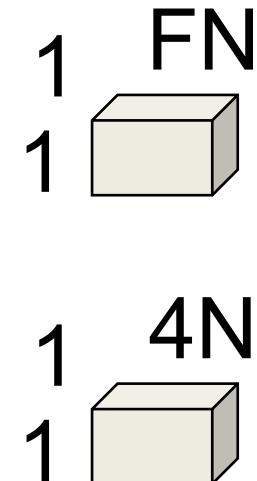
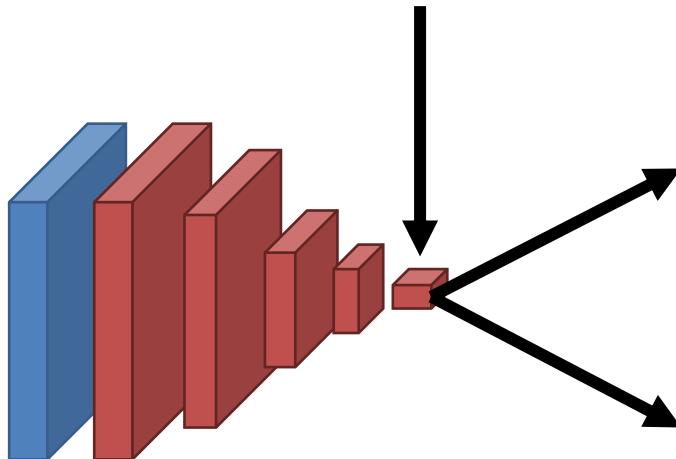


Now it's a herd of cows.
We need *lots* of outputs
(in fact the precise number of objects that are
in the image, which is circular reasoning).

In General

- Usually can't do varying-size outputs.
- Even if we could, think about how *you* would solve it if you were a network.

Bottleneck has to *encode* where the objects are for all objects and all N



An Alternate Approach

Examine every sub-window and determine if it is a tight box around an object



Yes



No?
Hold this thought



No

Sliding Window Classification

Let's assume we're looking for pedestrians in a box with a fixed aspect ratio.



Sliding Window

Key idea – just try all the subwindows in the image at all positions.



Generating hypotheses

Key idea – just try all the subwindows in the image at all positions **and scales**.



Note – Template did not change size

Each window classified separately



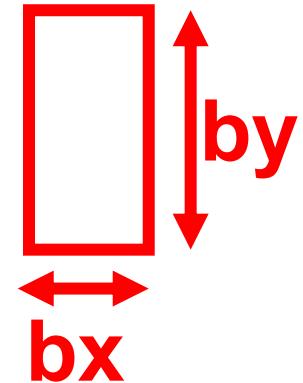
Slide credit: J. Hays

How Many Boxes Are There?

Given a $H \times W$ image and a “template” of size by, bx .

Q. How many sub-boxes are there of size (by, bx) ?

A. $(H-by)^*(W-bx)$

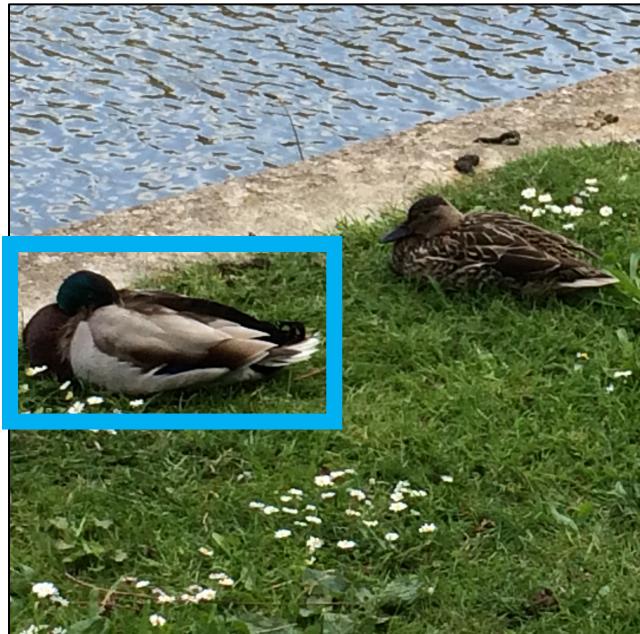


This is before considering adding:

- scales (by^*s, bx^*s)
- aspect ratios (by^*sy, bx^*sx)

Challenges of Object Detection

- Have to evaluate *tons* of boxes
- Positive instances of objects are *extremely* rare

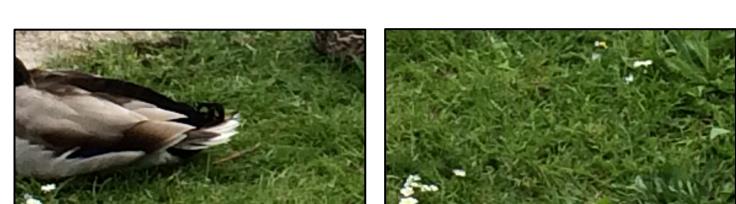
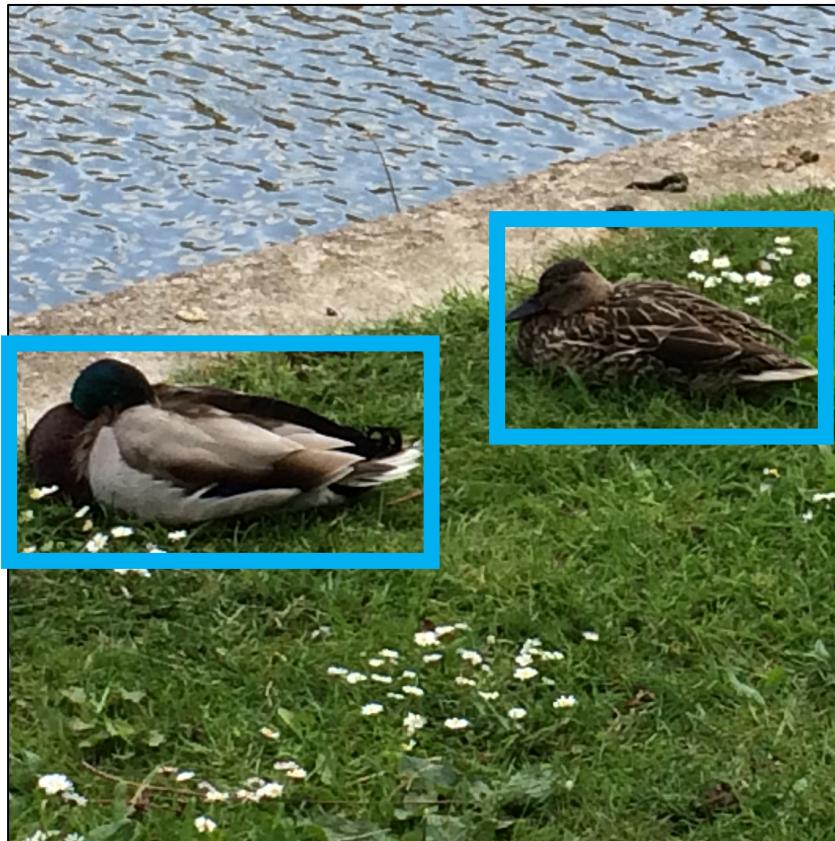


How many ways can we get the box wrong?

1. Wrong left x
2. Wrong right x
3. Wrong top y
4. Wrong bottom y

Evaluating – Bounding Boxes

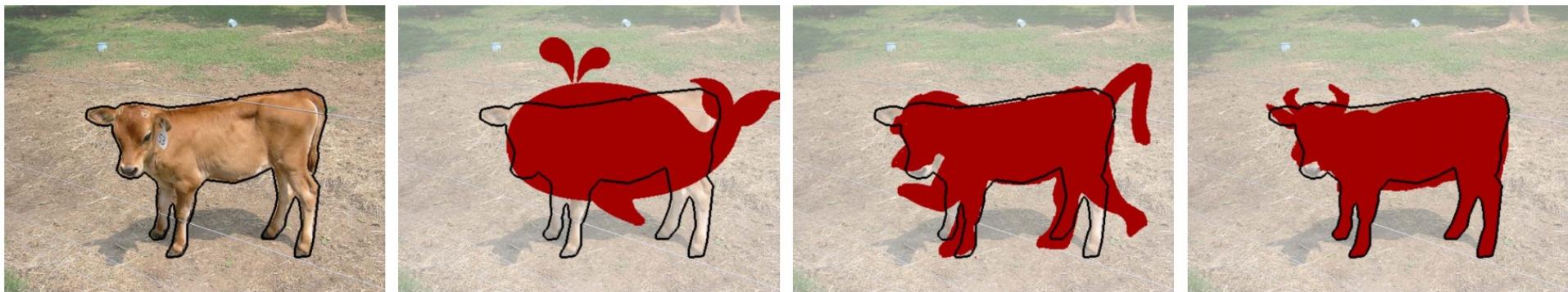
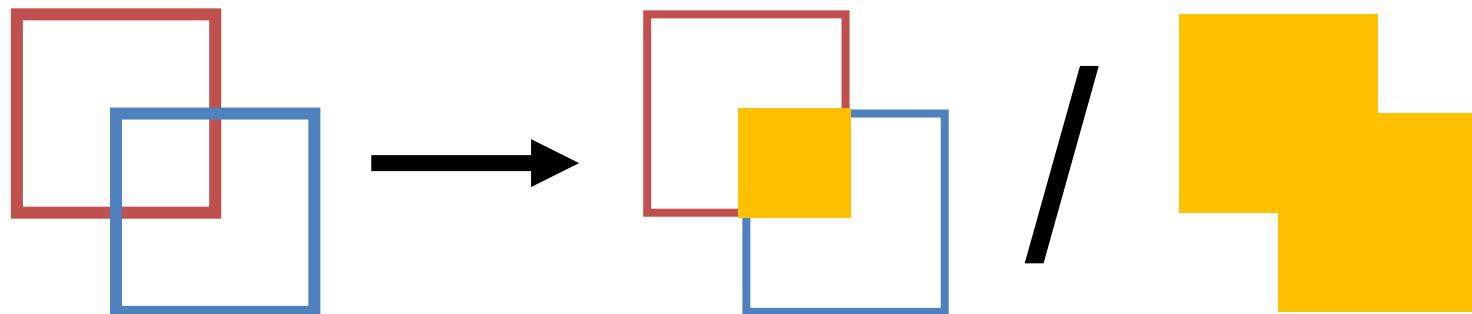
Raise your hand when you think the detection stops being correct.



Evaluating – Bounding Boxes

Standard metric for two boxes:

Intersection over union/IoU/Jaccard index



(a) Ground truth

(b) $\mathcal{J} = 0.554$

(c) $\mathcal{J} = 0.703$

(d) $\mathcal{J} = 0.910$

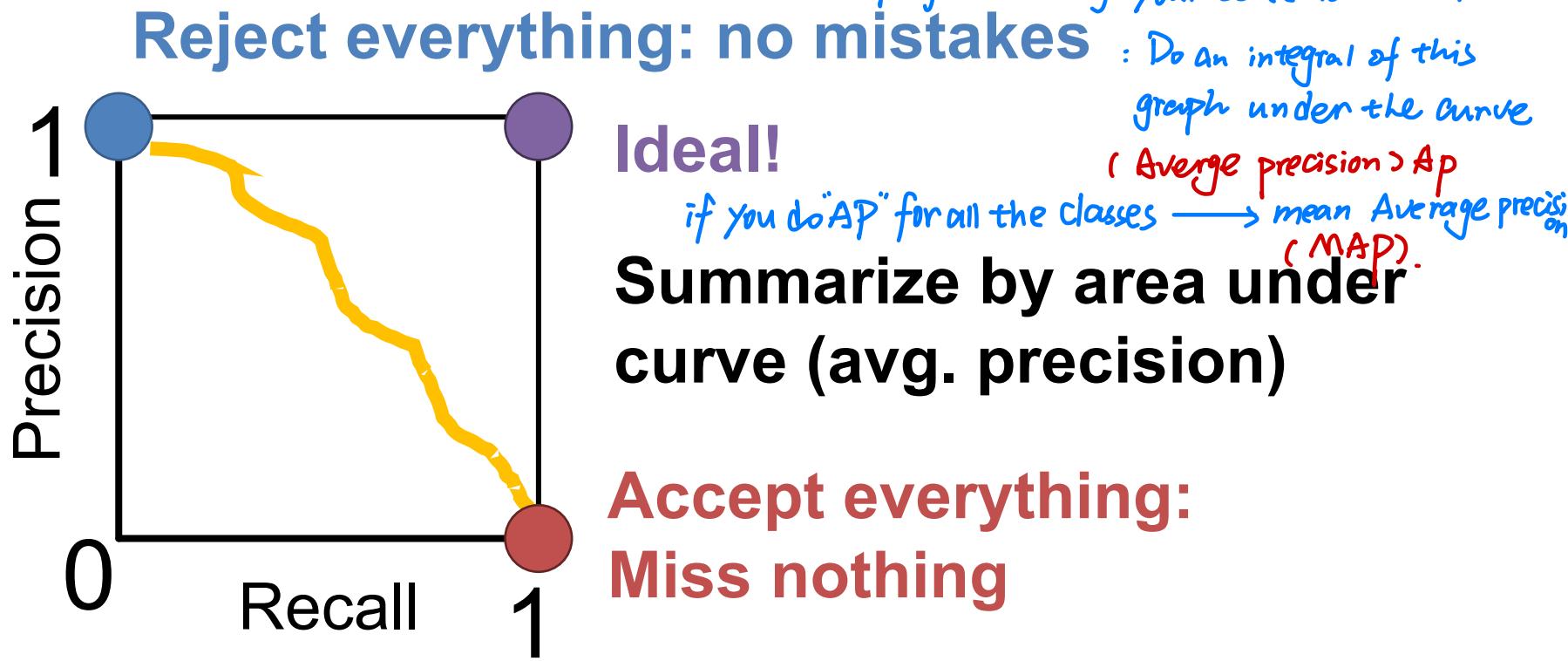
Evaluating Performance

- Remember: accuracy = average of whether prediction is correct
- Suppose I have a system that gets 99% accuracy in person detection.
- **What's wrong?**
- I can get that by just saying no object everywhere!

Evaluating Performance

- True detection aka true positive: high IoU (>0.5)
- Precision: #true detections / #detections by detector
- Recall: #true detections / #ground truth positives

*if you want to use one number to describe
the performance of your detection model.*



Generic object detection



Histograms of oriented gradients (HOG)

Partition image into blocks and compute histogram of gradient orientations in each block

$H \times W \times 3$ Image



$H' \times W' \times C'$ Image

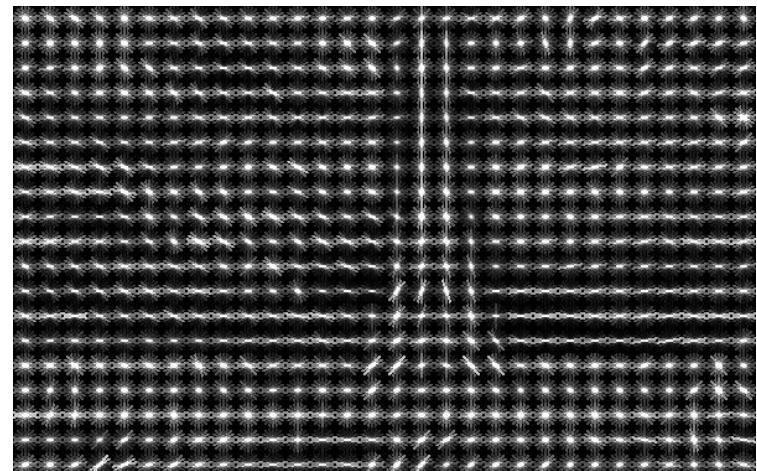


Image credit: N. Snavely

N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#),
CVPR 2005

Slide Credit: S. Lazebnik

Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine

positive training examples



negative training examples



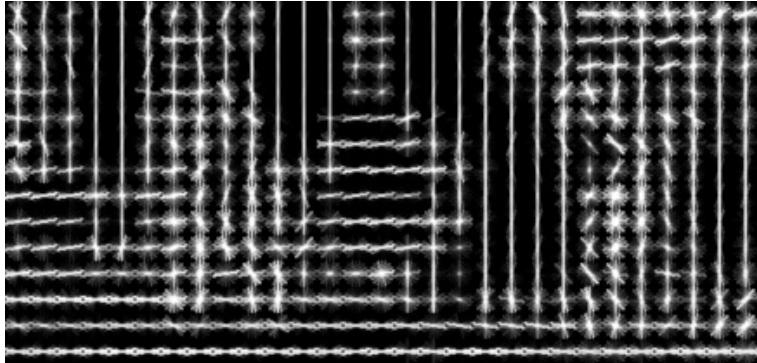
N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#),
CVPR 2005

Slide Credit: S. Lazebnik

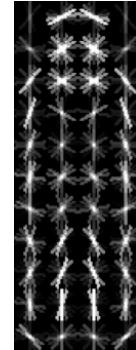
Pedestrian detection with HOG

- Train pedestrian “template” using a linear svm
- At test time, convolve feature map with template
- Find local maxima of response
- For multi-scale detection, repeat over multiple levels of a HOG *pyramid*

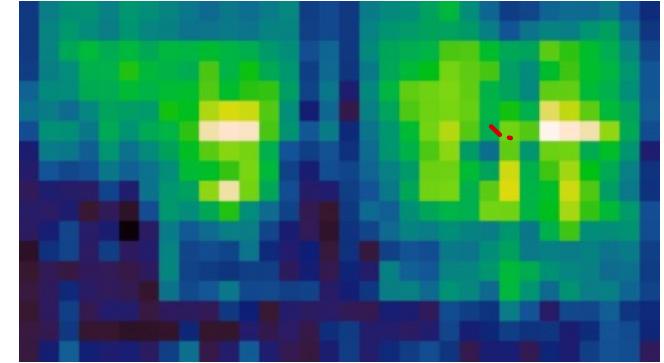
HOG feature map



Template



Detector response map



Example detections



[Dalal and Triggs, CVPR 2005]

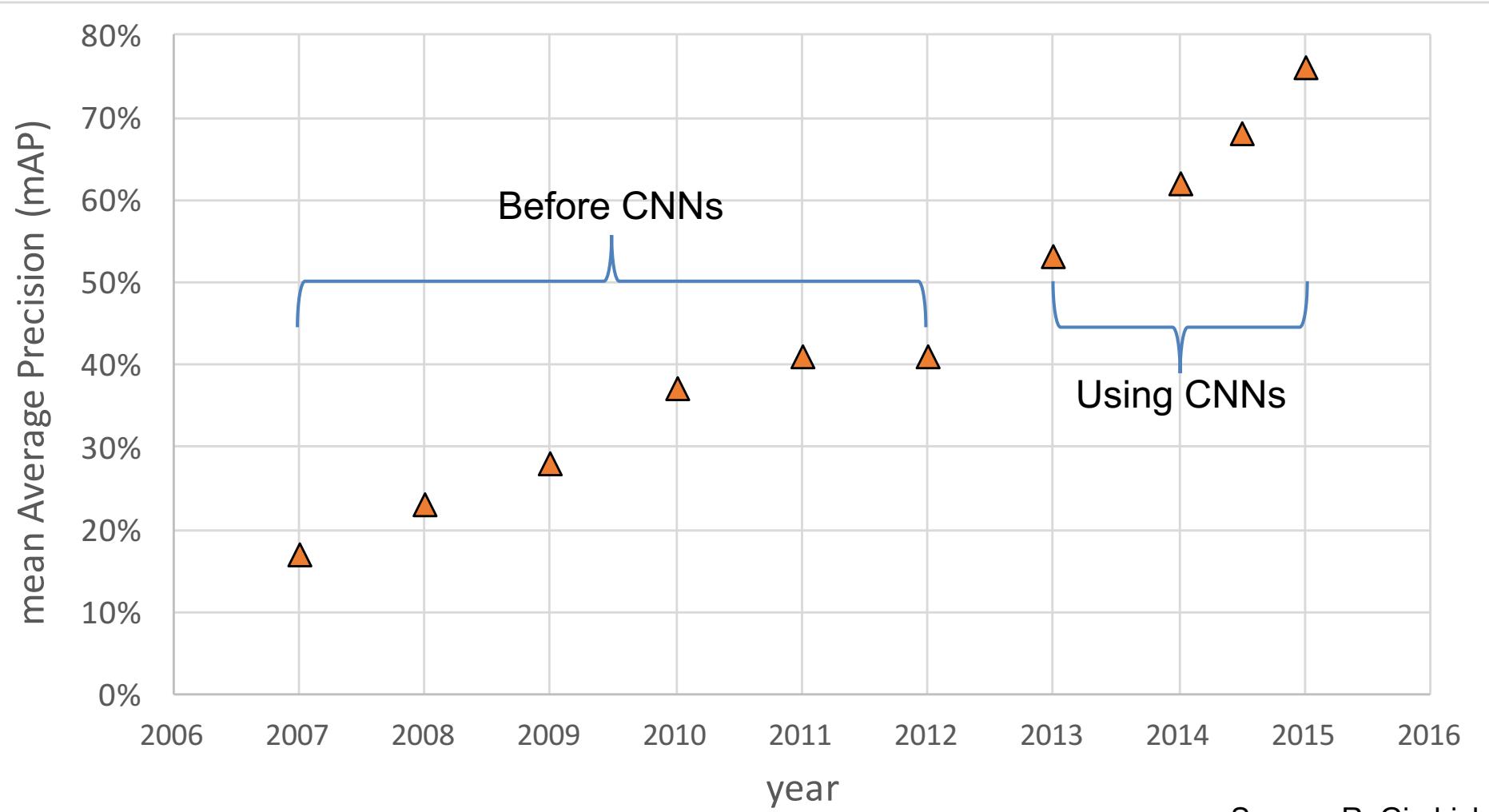
PASCAL VOC Challenge (2005-2012)



- 20 challenge classes:
- *Person*
- *Animals*: bird, cat, cow, dog, horse, sheep
- *Vehicles*: aeroplane, bicycle, boat, bus, car, motorbike, train
- *Indoor*: bottle, chair, dining table, potted plant, sofa, tv/monitor
- Dataset size (by 2012): 11.5K training/validation images, 27K bounding boxes, 7K segmentations

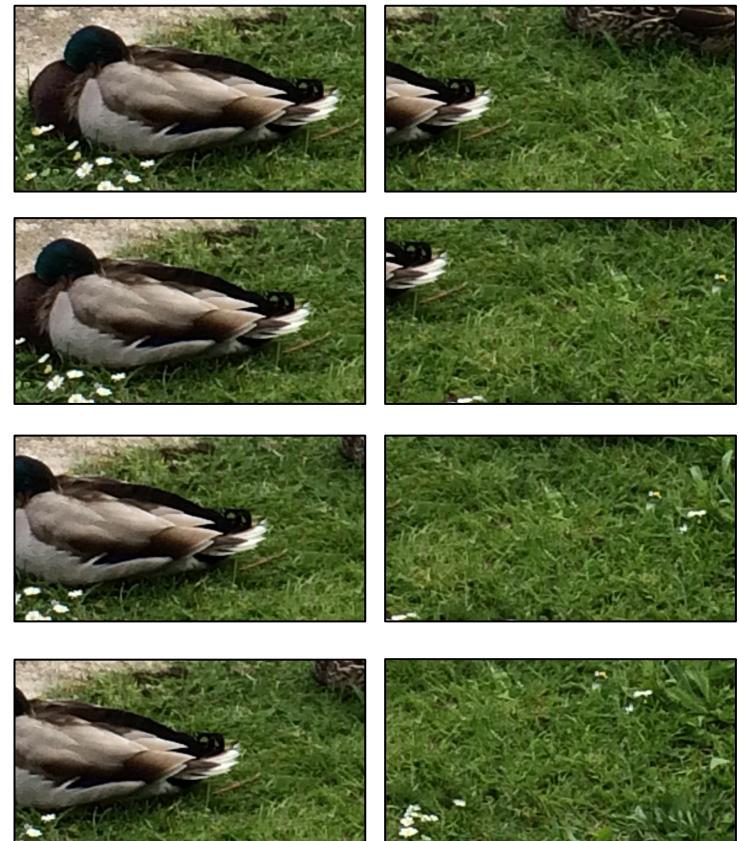
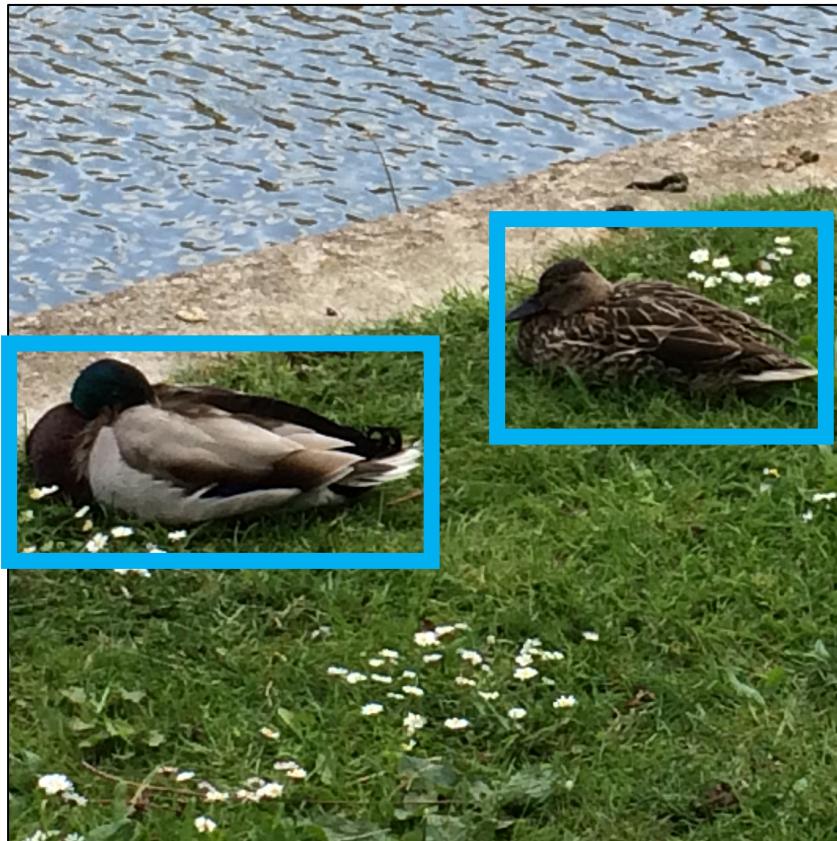
Object detection progress

PASCAL VOC

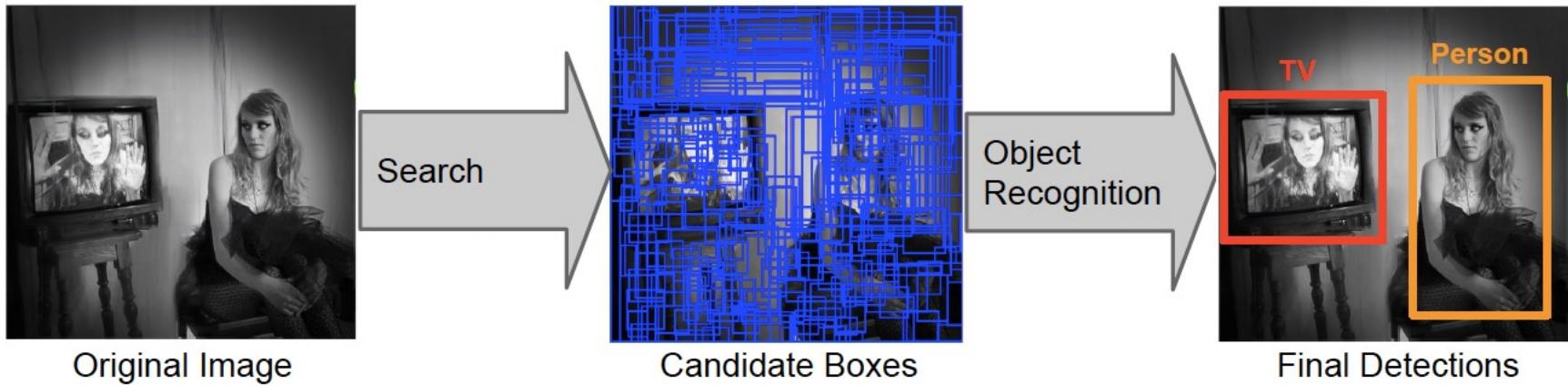


Region Proposals

Do I need to spend a lot of time filtering all the boxes covering grass?



Region Proposals



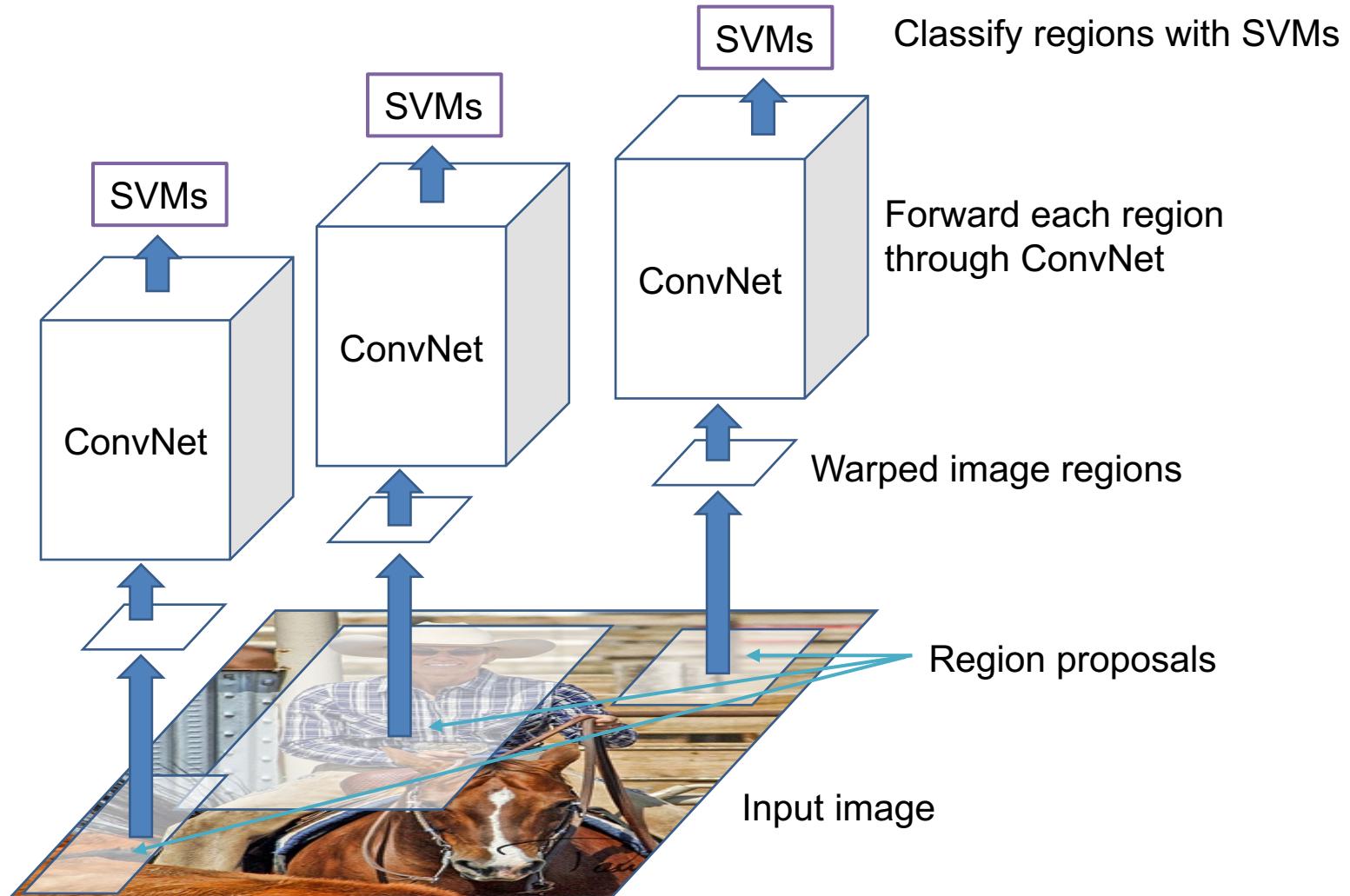
- As an alternative to sliding window search, evaluate a few hundred *region proposals*
 - Can use slower but more powerful features and classifiers
 - Proposal mechanism can be category-independent
 - Proposal mechanism can be trained

Slide Credit: S. Lazebnik

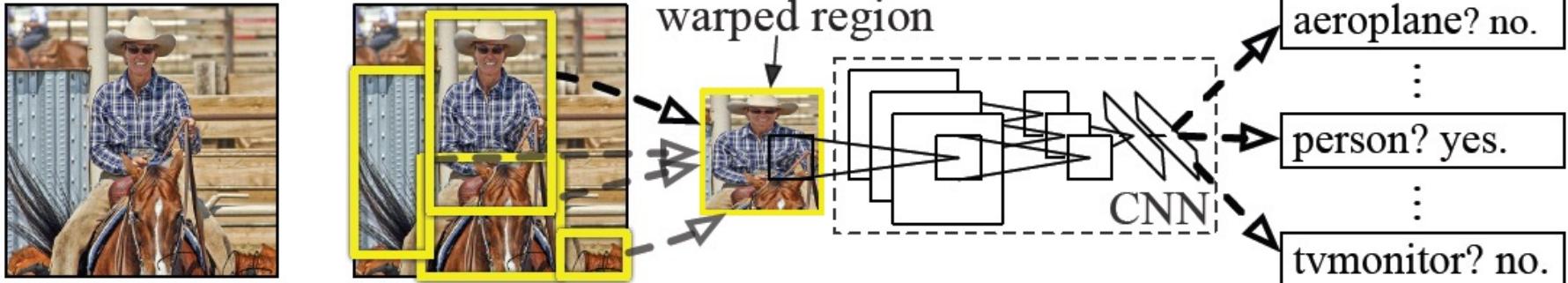
Example Region Proposal Algo: Edge boxes: Locating object proposals from edges. 2014

R-CNN: Region proposals + CNN features

Source: R. Girshick



R-CNN details

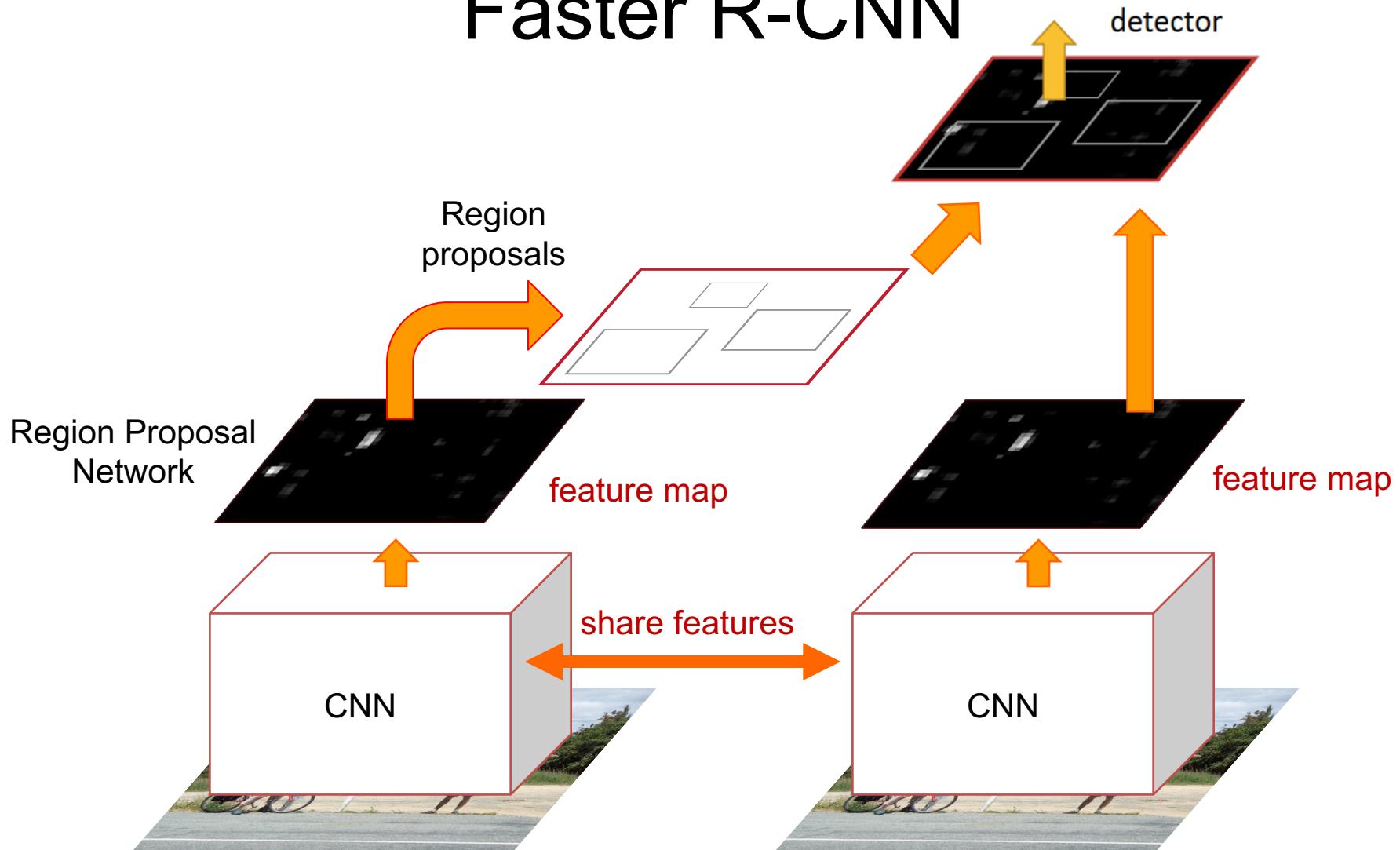


- **Regions:** ~2000 Selective Search proposals
- **Network:** AlexNet *pre-trained* on ImageNet (1000 classes), *fine-tuned* on PASCAL (21 classes)
- **Final detector:** warp proposal regions, extract fc7 network activations (4096 dimensions), classify with linear SVM
- **Bounding box regression** to refine box locations
- **Performance:** mAP of 53.7% on PASCAL 2010 (vs. 35.1% for Selective Search and 33.4% for DPM).

R-CNN pros and cons

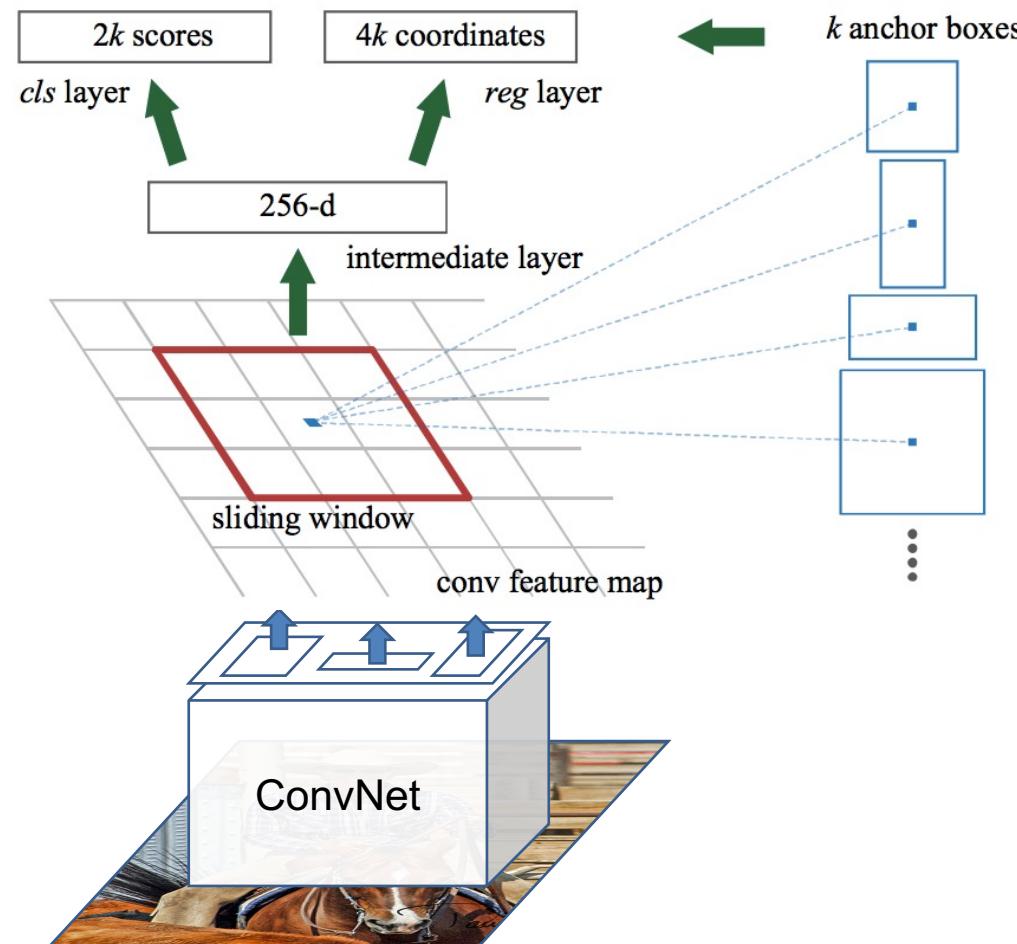
- Pros
 - Accurate!
 - Any deep architecture can immediately be “plugged in”
- Cons
 - Ad hoc training objectives
 - Fine-tune network with softmax classifier (log loss)
 - Train post-hoc linear SVMs (hinge loss)
 - Train post-hoc bounding-box regressions (least squares)
 - Training is slow (84h), takes a lot of disk space
 - 2000 CNN passes per image
 - Inference (detection) is slow (47s / image with VGG16)

Faster R-CNN



S. Ren, K. He, R. Girshick, and J. Sun, [Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks](#), NIPS 2015

Region Proposal Network (RPN)

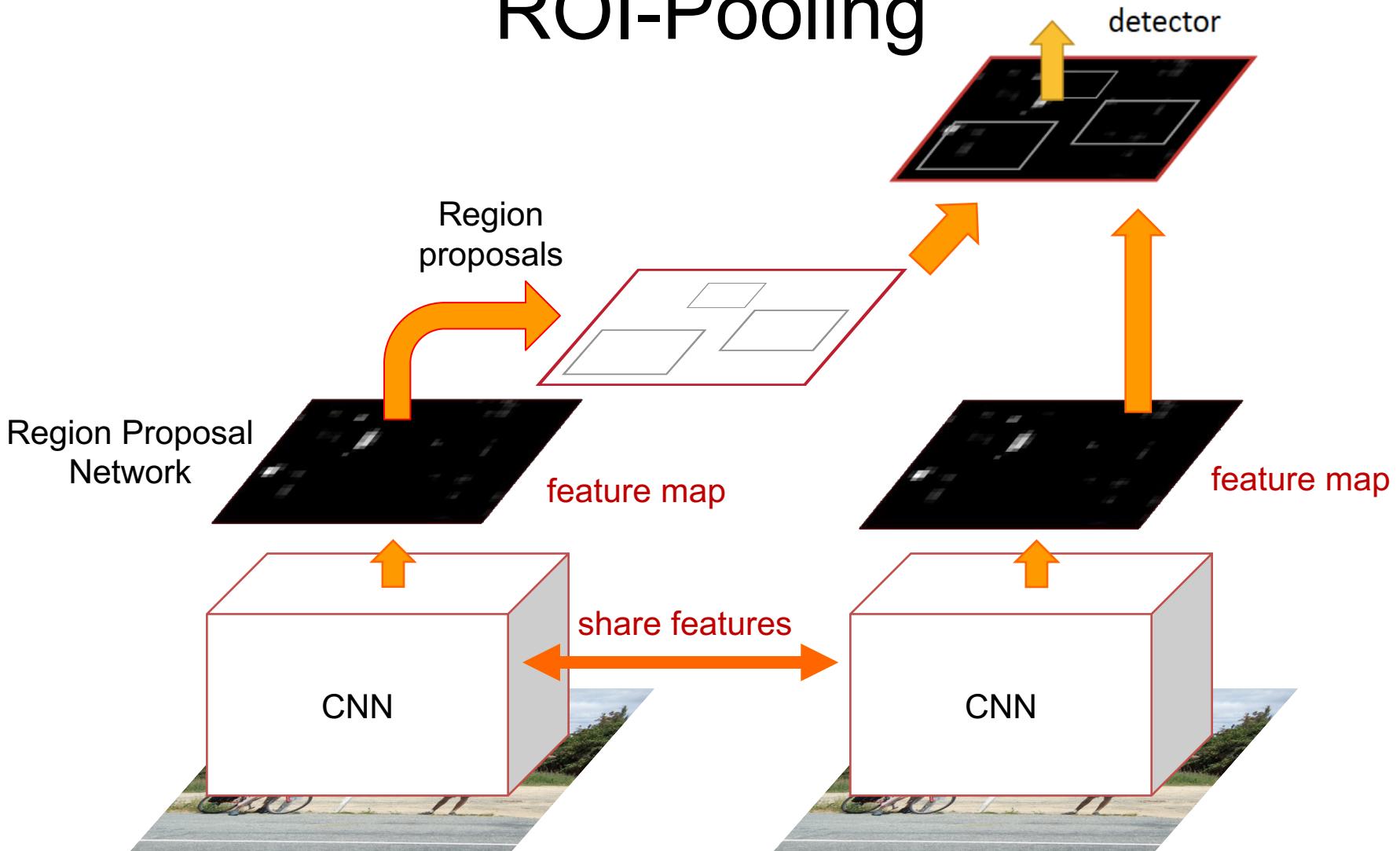


Small network applied to conv5 feature map.

Predicts:

- good box or not (classification),
 - how to modify box (regression)
- for k “anchors” or boxes relative to the position in feature map.

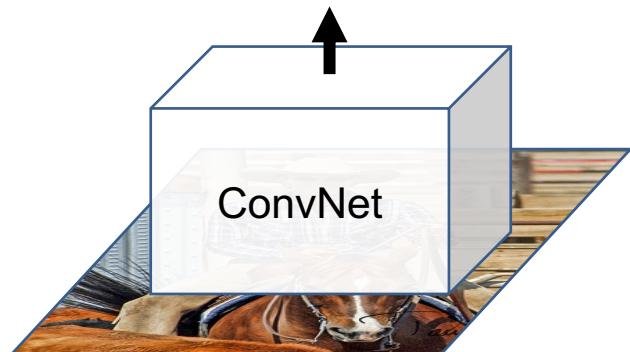
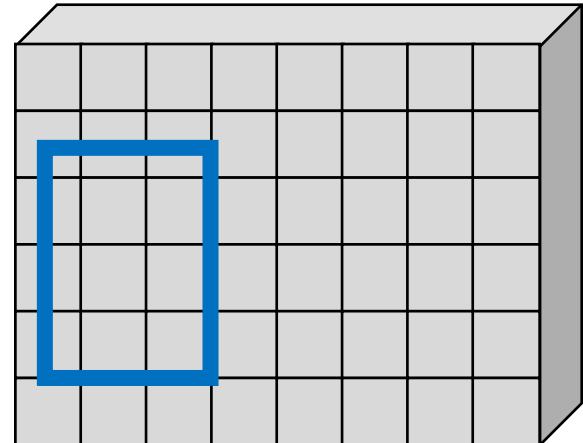
ROI-Pooling



S. Ren, K. He, R. Girshick, and J. Sun, [Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks](#), NIPS 2015

ROI Pooling/Align

Feature Map
(e.g., 6x8x256)

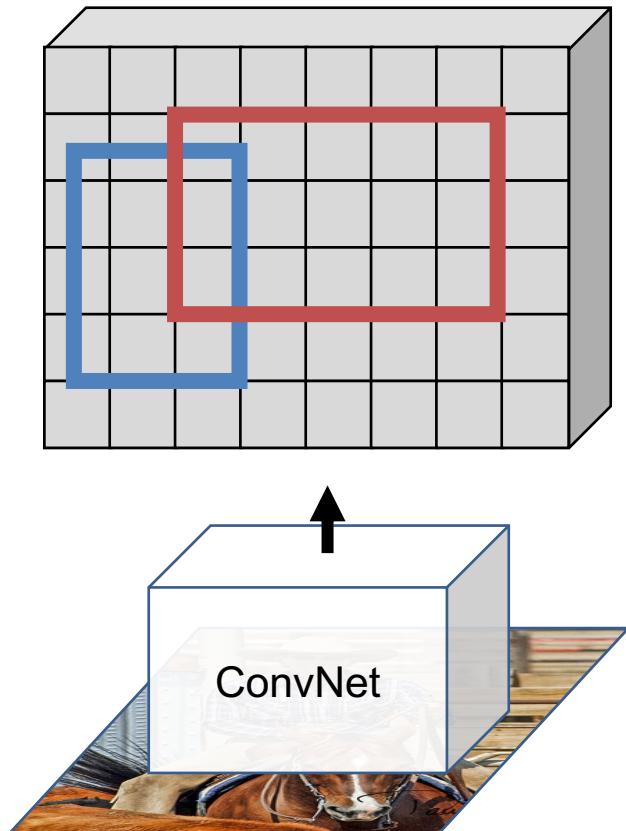


Given box in original image,
calculate where the box goes.

- Example: $H=600, W=800$
- Feature map is $H'=6, W'=8$
- Box: left $x=50$, top $y=150$,
 $\text{width}=250$, $\text{height}=350$
- Feature map box: left $x=0.5$, $y=1.5$, $\text{width}=2.5$, $\text{height}=3.5$

ROI Pooling/Align

Feature Map
(e.g., 6x8x256)

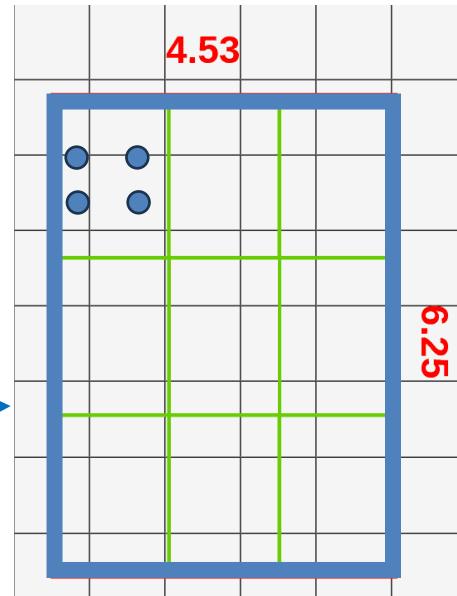
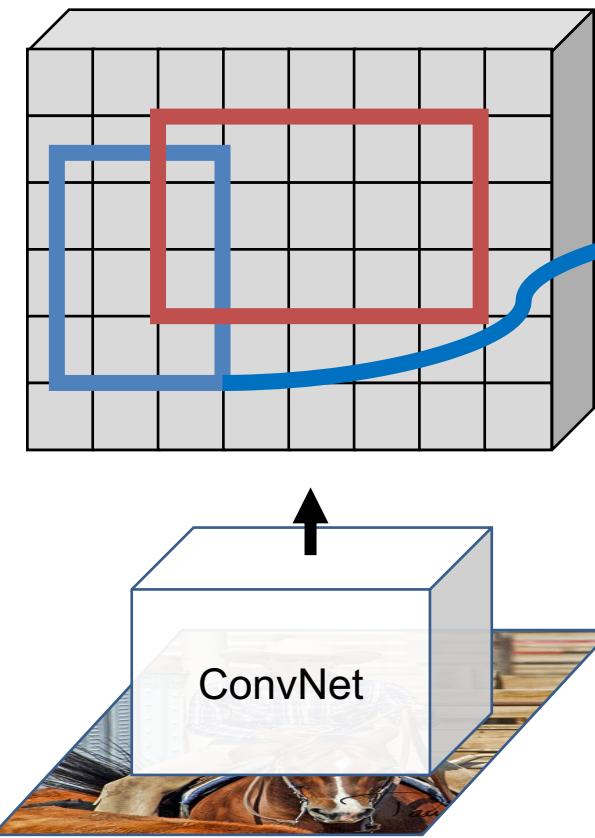


Given box in original image,
calculate where the box goes.

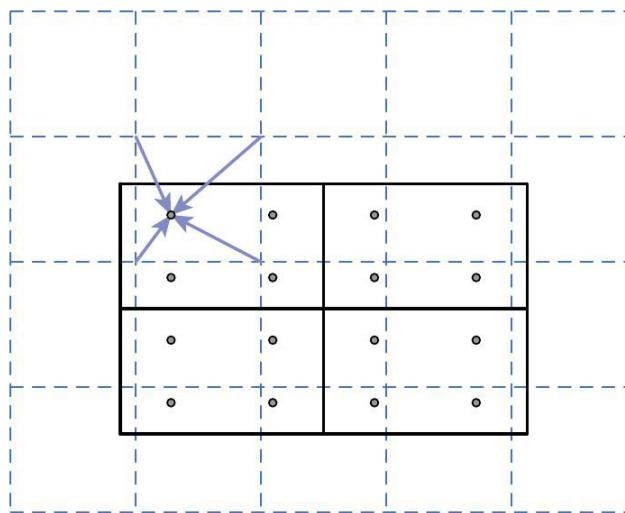
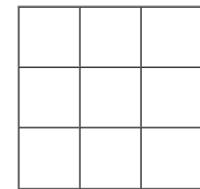
- Example: $H=600, W=800$
- Feature map is $H'=6, W'=8$
- Box: left $x=50$, top $y=150$,
width=250, height=350
- Feature map box: left $x=0.5$, $y = 1.5$, width=2.5, height=3.5
- Other feature map box: left $x=2$,
top $y = 1$, width=5, height=3

ROI Pooling/Align

Feature Map

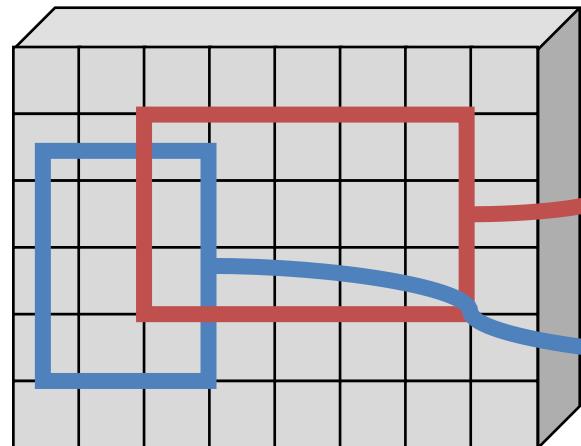


3x3 Pooled Feature

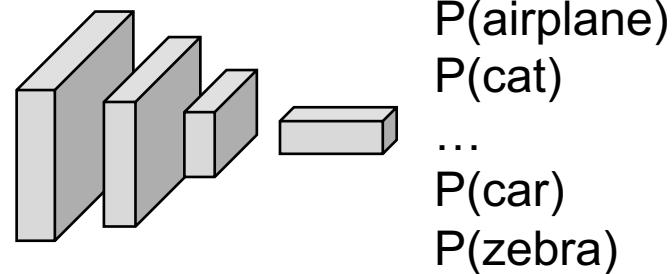


ROI Pooling/Align

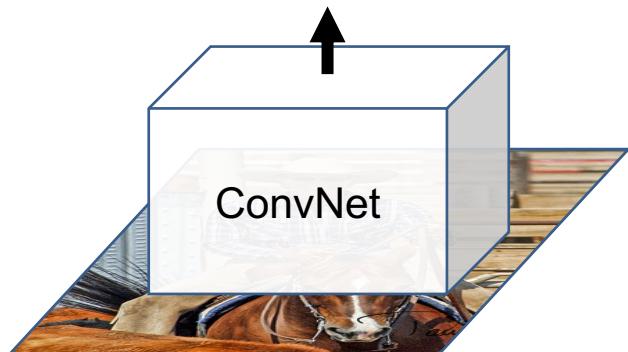
Feature Map
(e.g., 6x8x256)



Resize to fixed size (e.g., 7x7)
Details critical, but beyond scope of class.

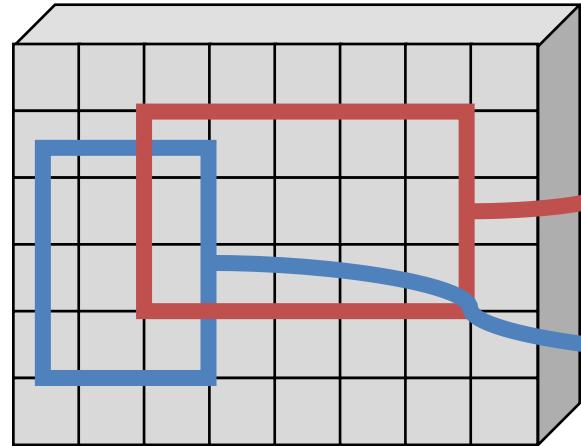


Afterwards, can add a small neural network that classifies the box and is applied to each window.

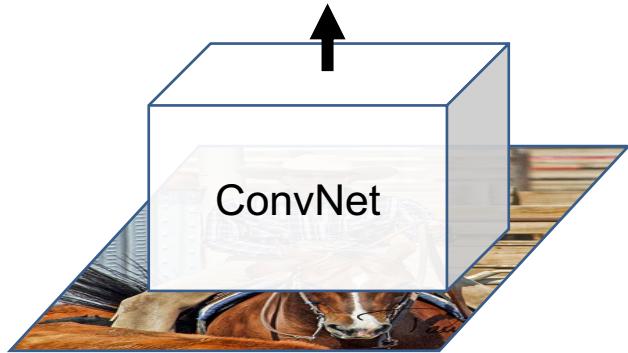
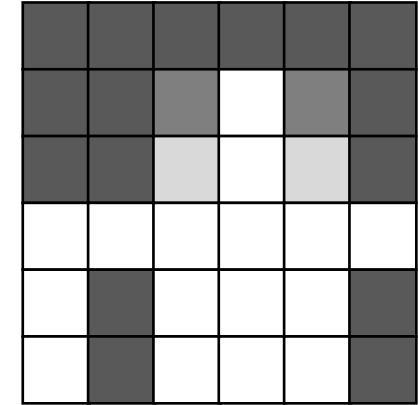
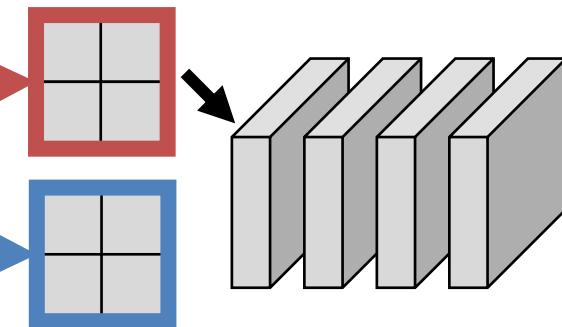


Mask RCNN

Feature Map
(e.g., 6x8x256)



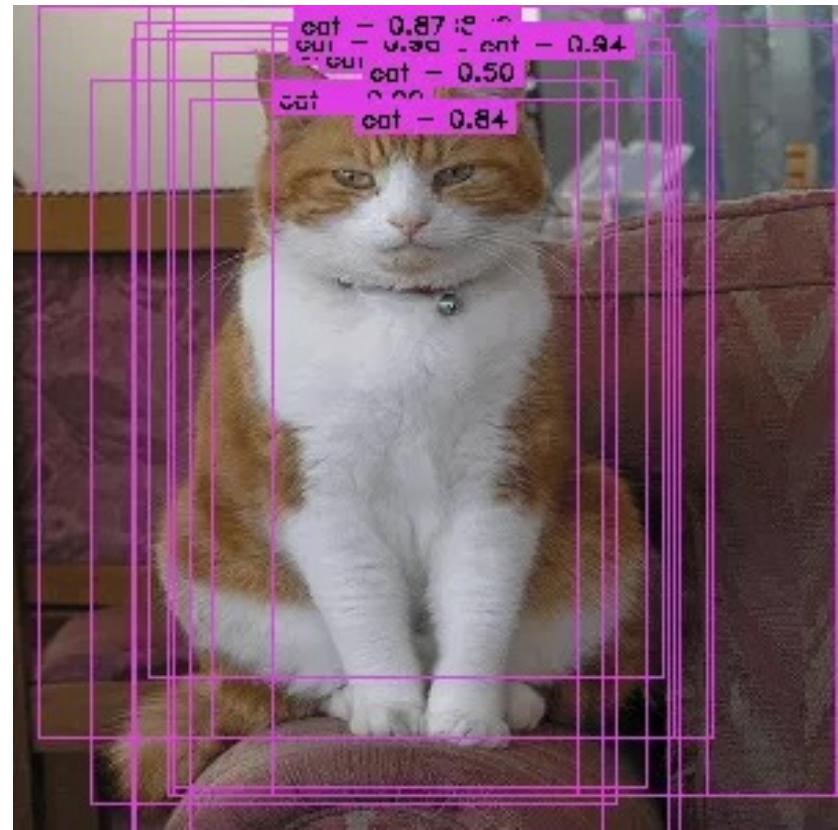
Resize to fixed size (e.g., 7x7)
Details critical, but beyond
scope of class.



Can also predict a mask!
Everything is learned together
Simple and effective; details critical

Non-Maximum Supression

1. Select ROI with highest confidence
2. Discard all ROIs with significant overlap (IoU)
3. Go to 1.

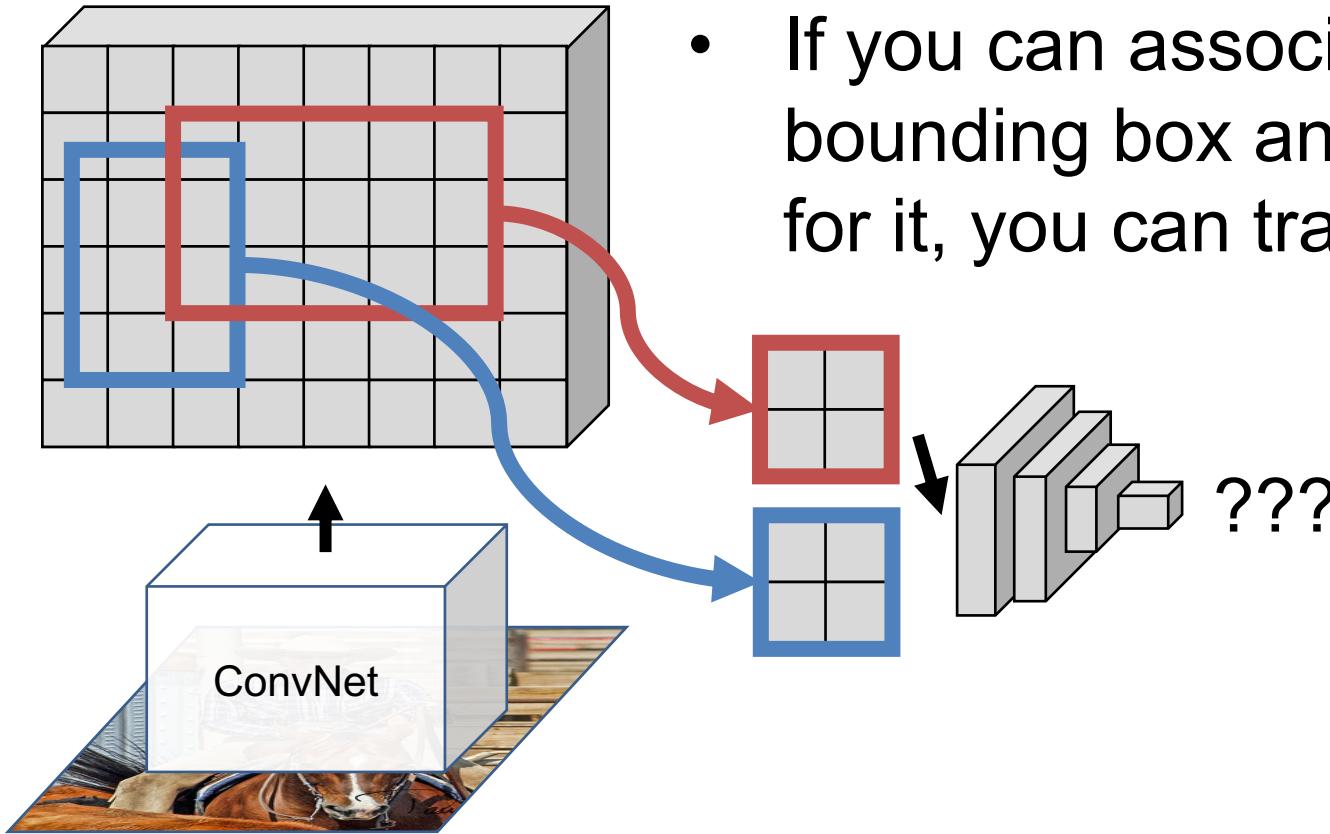


MaskRCNN – Results



Extending Object Detection

Feature Map
(e.g., 6x8x256)



- Can ask the network to predict *nearly* anything.
- If you can associate it with a bounding box and can get data for it, you can train the model

Extending Object Detection

Example: RGB image input, detect planar surfaces



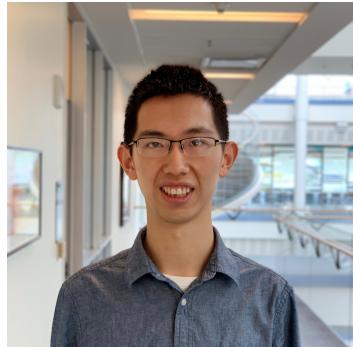
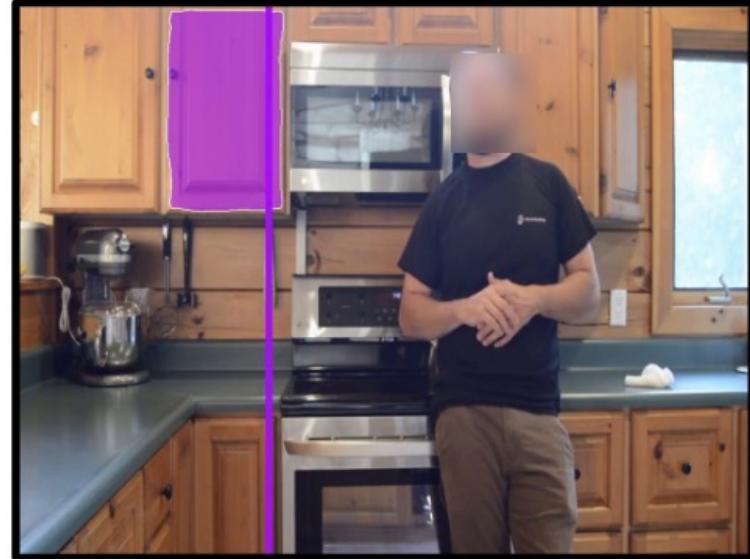
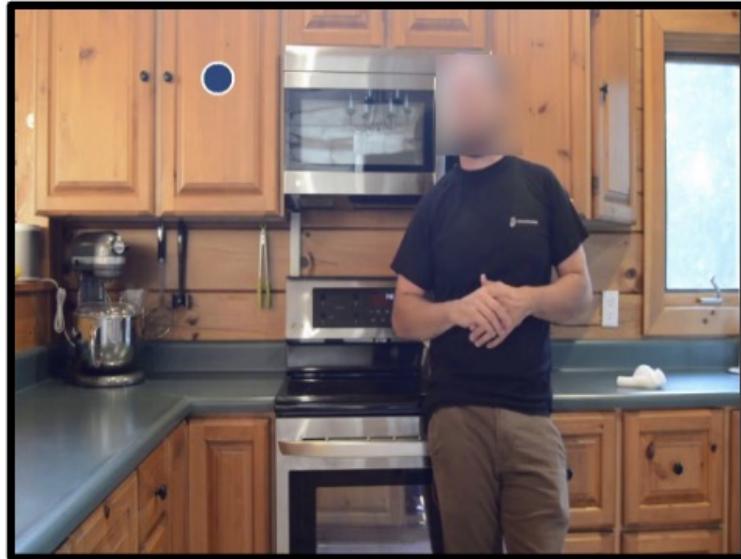
Extending Object Detection



Core building
block is detecting
plane in image.

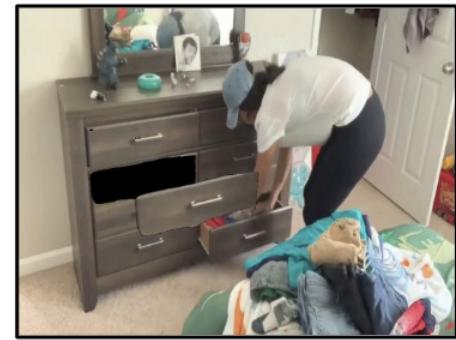
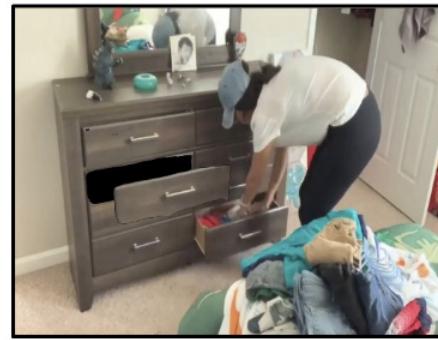
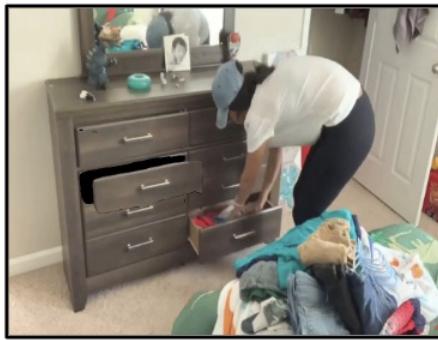
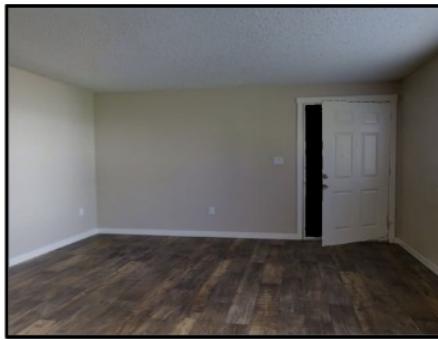


Other Fun Stuff

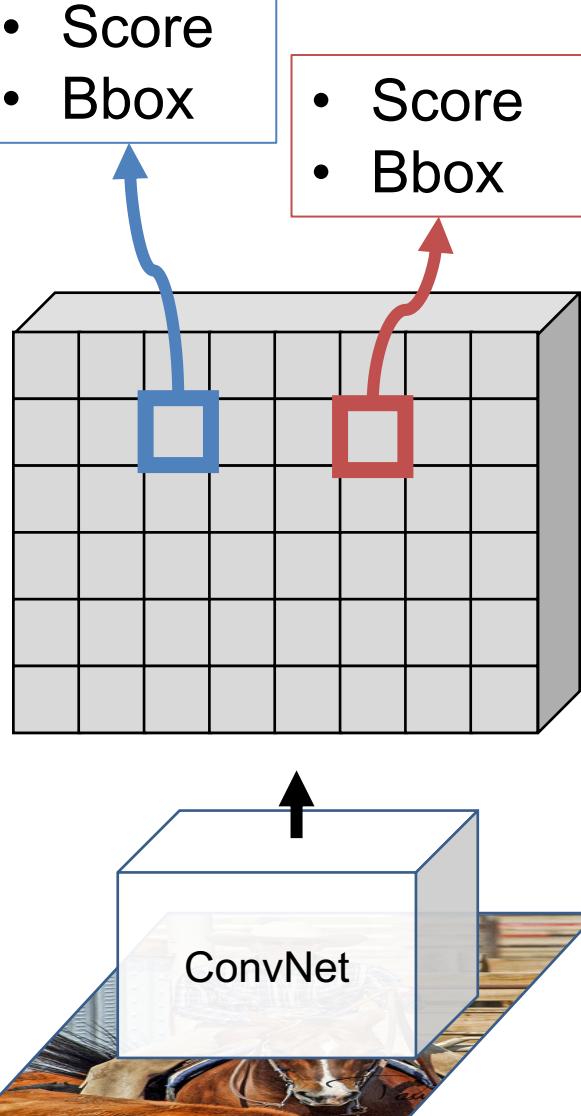


Shengyi Qian
Former 442 IA!

Other Fun Stuff



YOLO

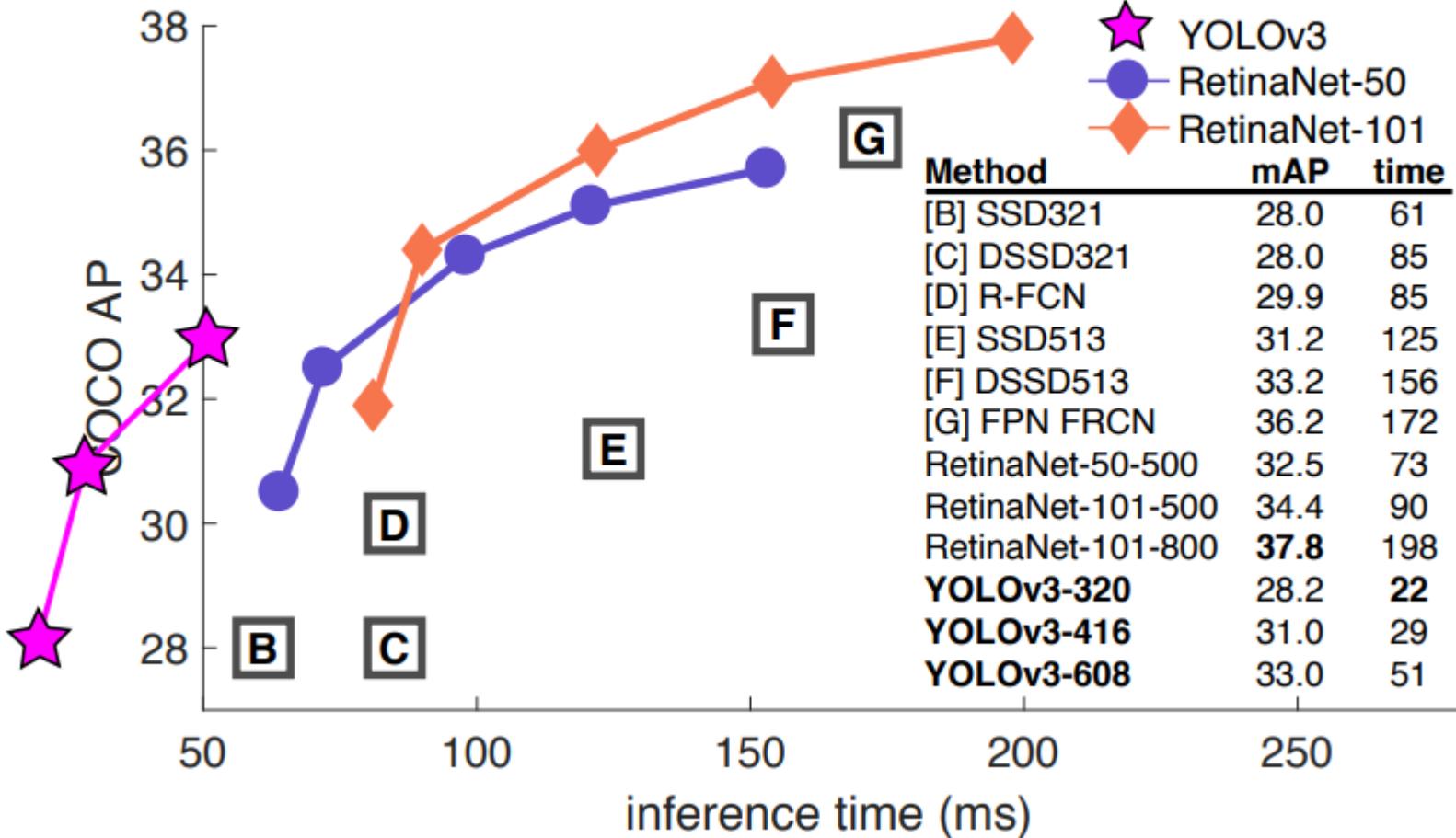


- No proposals
- Predict at each location in 7×7 feature map, score for each class + 2 bboxes
- 7x faster than Faster-RCNN, but worse accuracy, precision
- Immensely popular in robotics
- Loads of similar methods (YOLOv2, YOLOv3)

YOLO

- Score

7



1
S

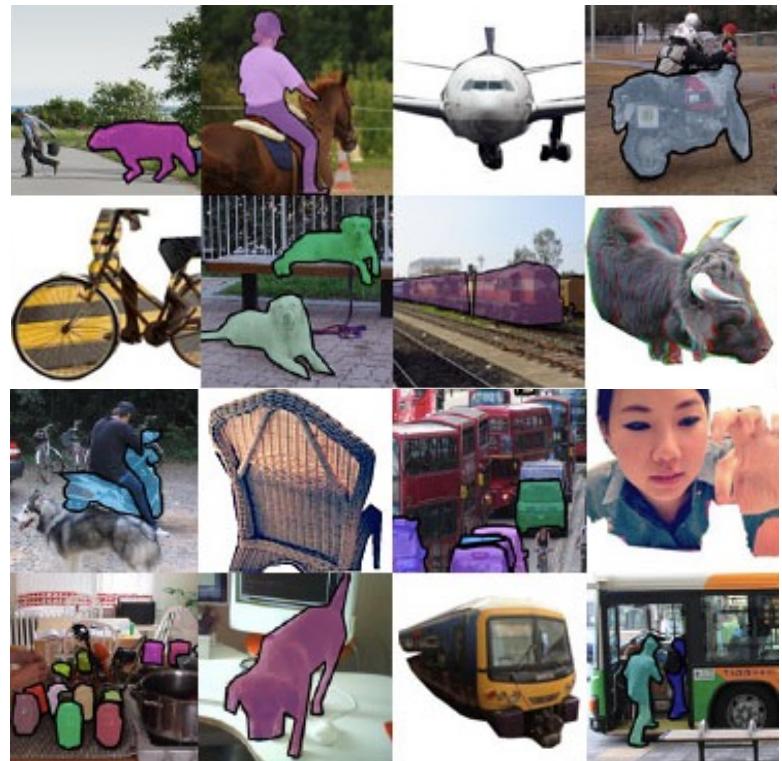
New detection benchmark: COCO (2014)

- 80 categories instead of PASCAL's 20
- Current best mAP: 66%



COCO is a large-scale object detection, segmentation, and captioning dataset. COCO has several features:

- ✓ Object segmentation
- ✓ Recognition in context
- ✓ Superpixel stuff segmentation
- ✓ 330K images (>200K labeled)
- ✓ 1.5 million object instances
- ✓ 80 object categories
- ✓ 91 stuff categories
- ✓ 5 captions per image
- ✓ 250,000 people with keypoints



<http://cocodataset.org/#home>

A Few Caveats

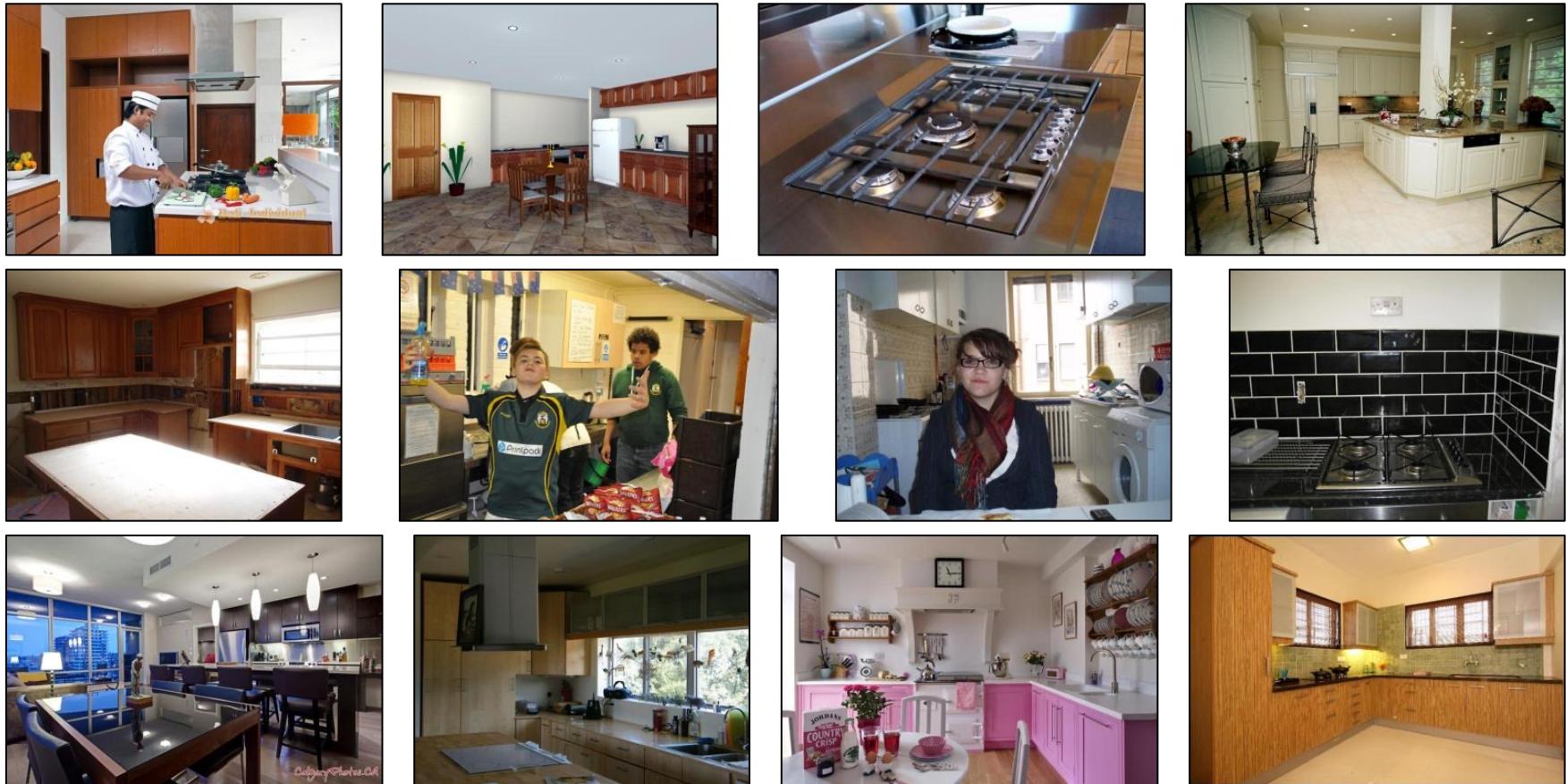
- Flickr images come from a really weird process
- Step 1: user takes a picture
- Step 2: user decides to upload it
- Step 3: user decides to write something like “refrigerator” somewhere in the description
- Step 4: a vision person stumbles on it while searching Flickr for refrigerators for a dataset



Who takes photos of open refrigerators ?????



Kitchens from Googling



Places 365 Dataset, Zhou et al. '17

Guess the category!

These were detected with >90% confidence



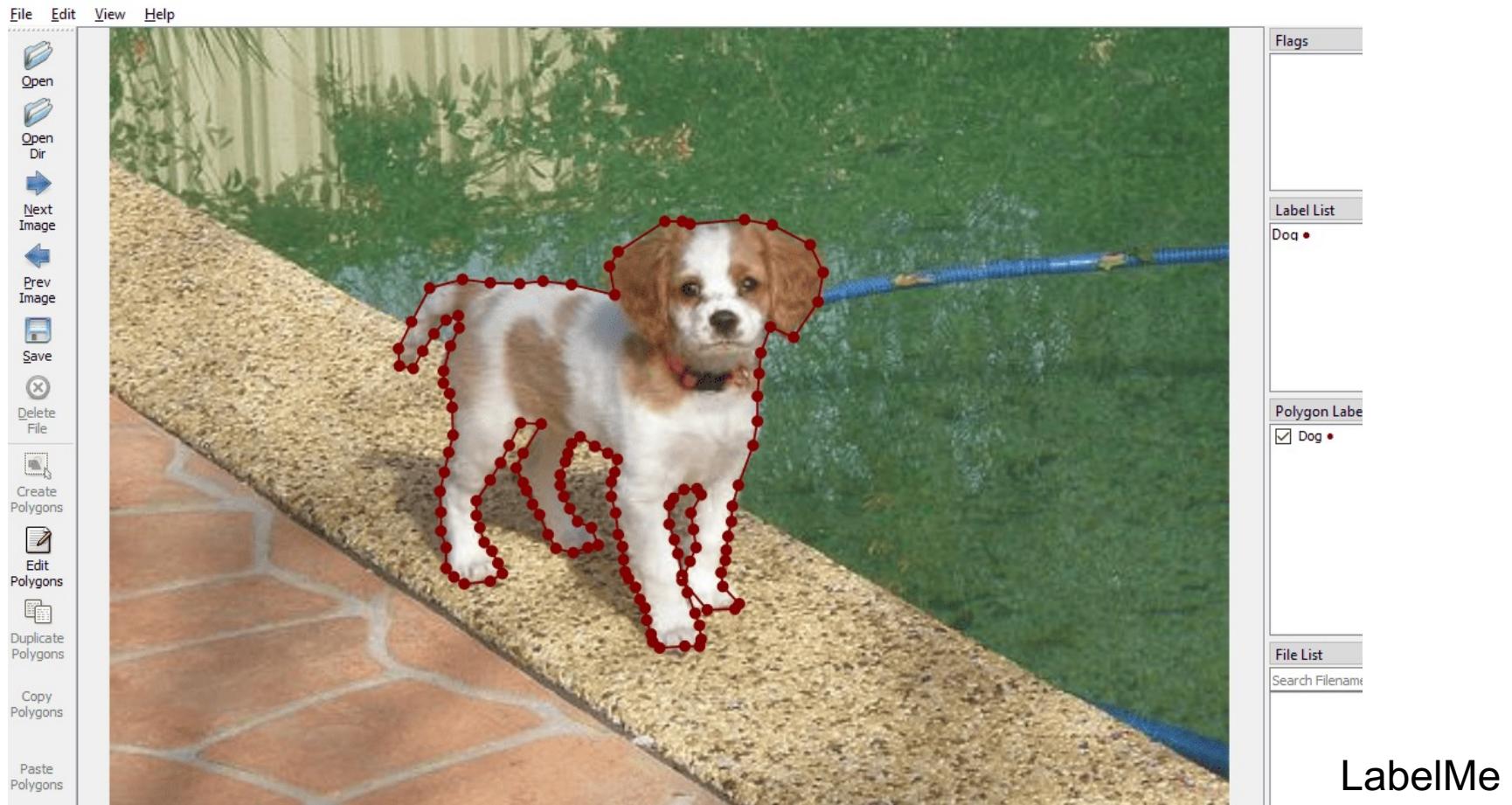
(1) Person

(2) Giraffe

(3) Bicycle

Self-Supervised Feature Learning

- Can we train without expensive labeling?



How Much Information is the Machine Given during Learning?

► “Pure” Reinforcement Learning (**cherry**)

- The machine predicts a scalar reward given once in a while.

► **A few bits for some samples**



► Supervised Learning (**icing**)

- The machine predicts a category or a few numbers for each input
- Predicting human-supplied data
- **10→10,000 bits per sample**

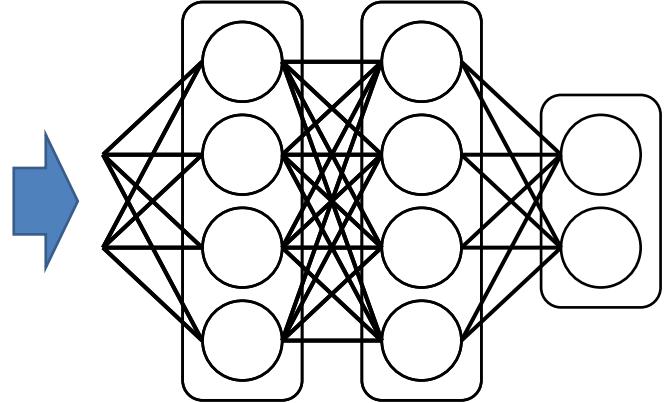
► Self-Supervised Learning (**cake génoise**)

- The machine predicts any part of its input for any observed part.
- Predicts future frames in videos
- **Millions of bits per sample**

Transfer Learning

- $w = \text{weights_from_somewhere_else}$
- for batch in batches:
 - inputs, labels = batch
 - calculate gradient of loss function with respect to w applied to samples in inputs
 - $w += \text{gradient}$

ImageNet + Deep Learning

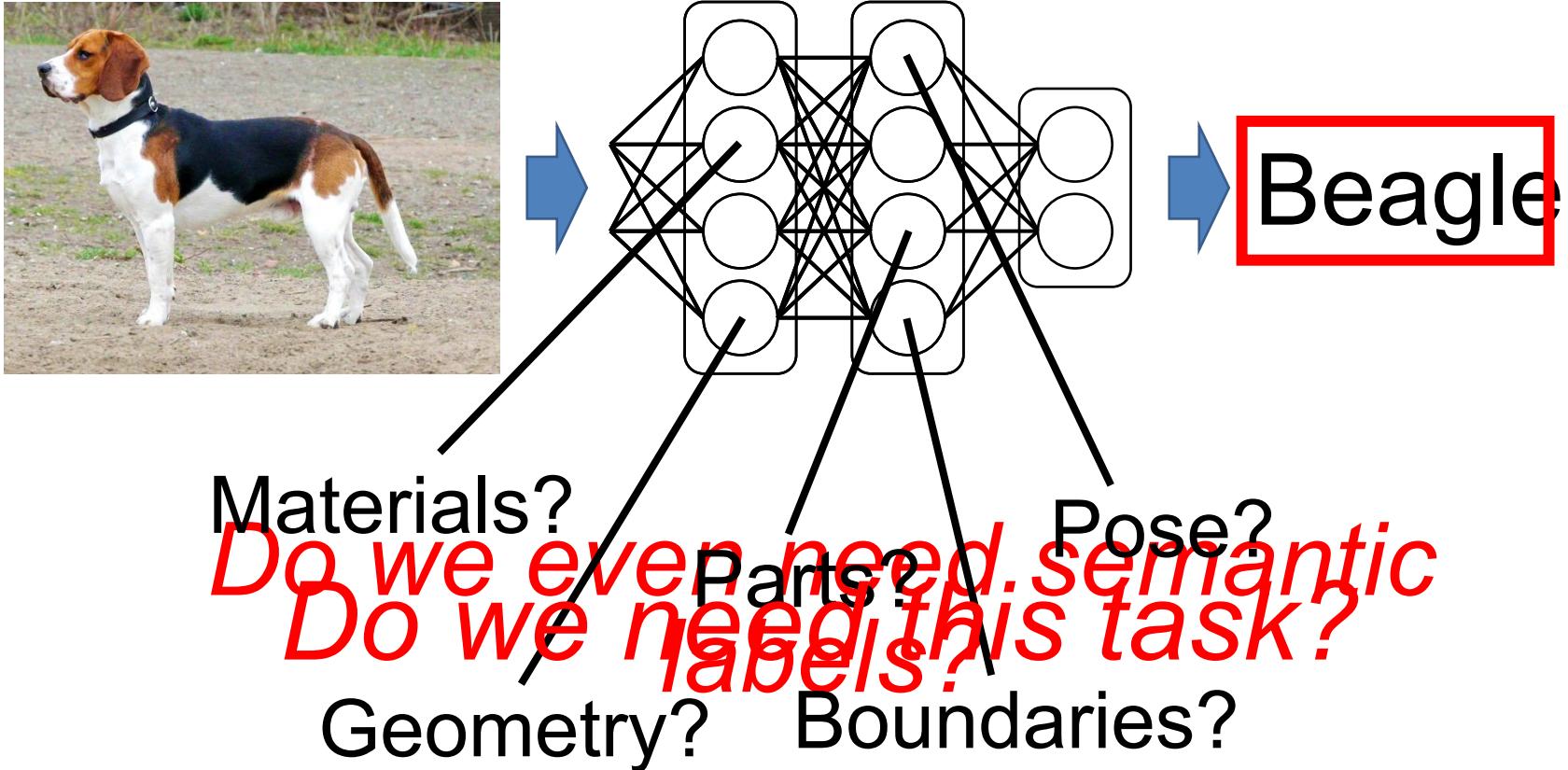


Beagle



- Image Retrieval
- Detection
- Segmentation
- Depth Estimation
- ...

ImageNet + Deep Learning



Context as Supervision

[Collobert & Weston 2008; Mikolov et al. 2013]

house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal milk; but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and appliances, the toasters and carpet

Deep
Net

Context Prediction for Images

?



?

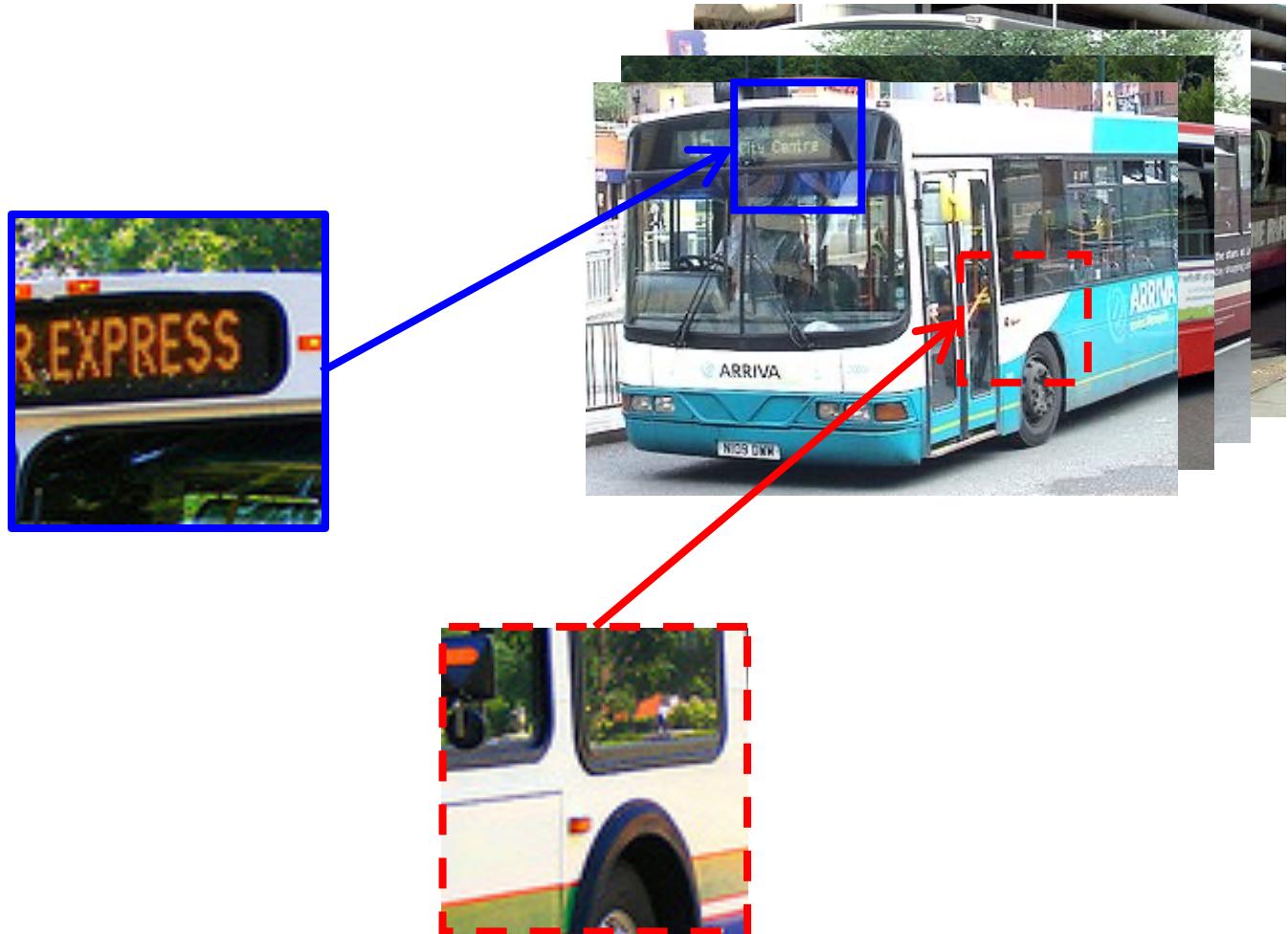
A

?

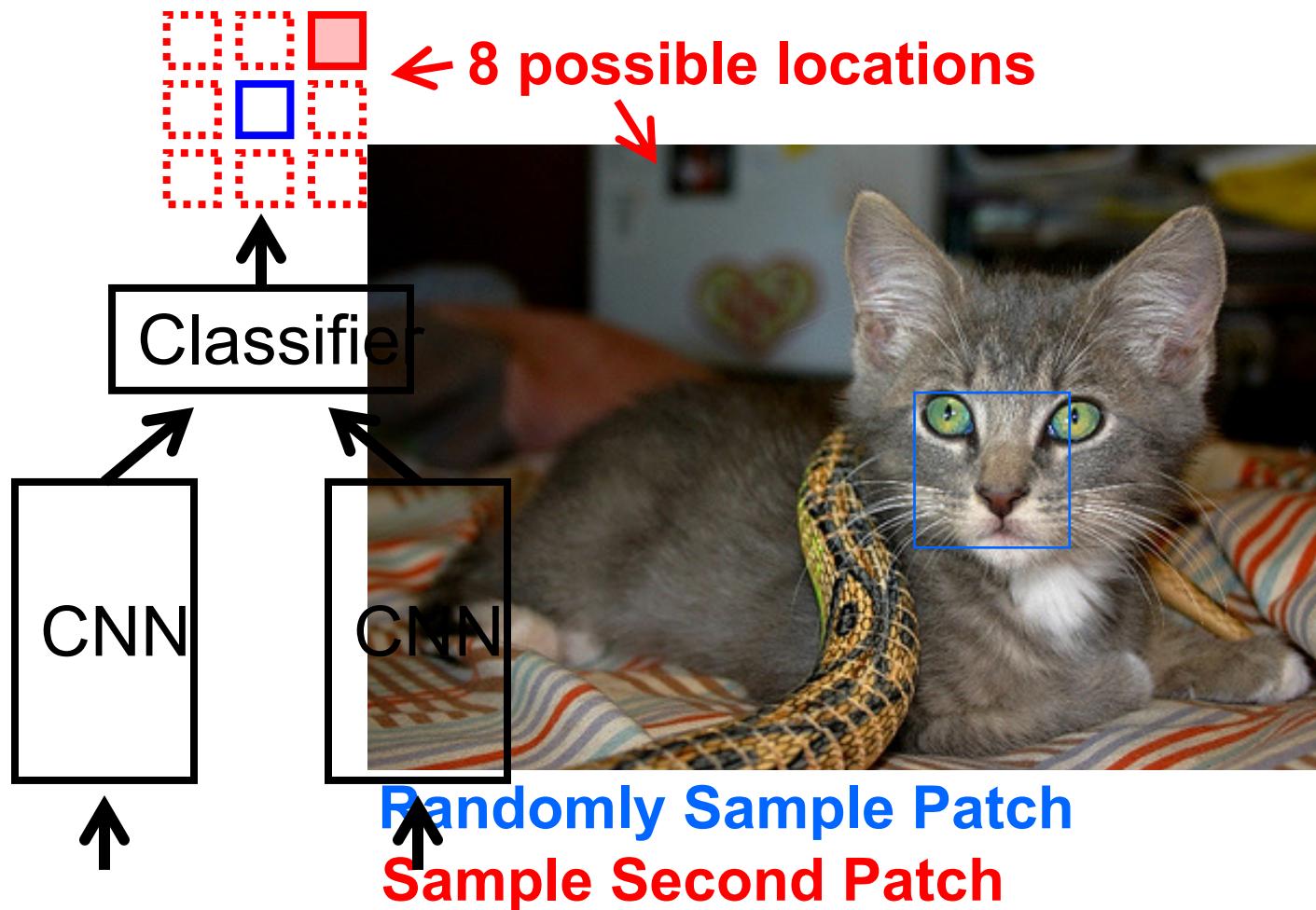
B

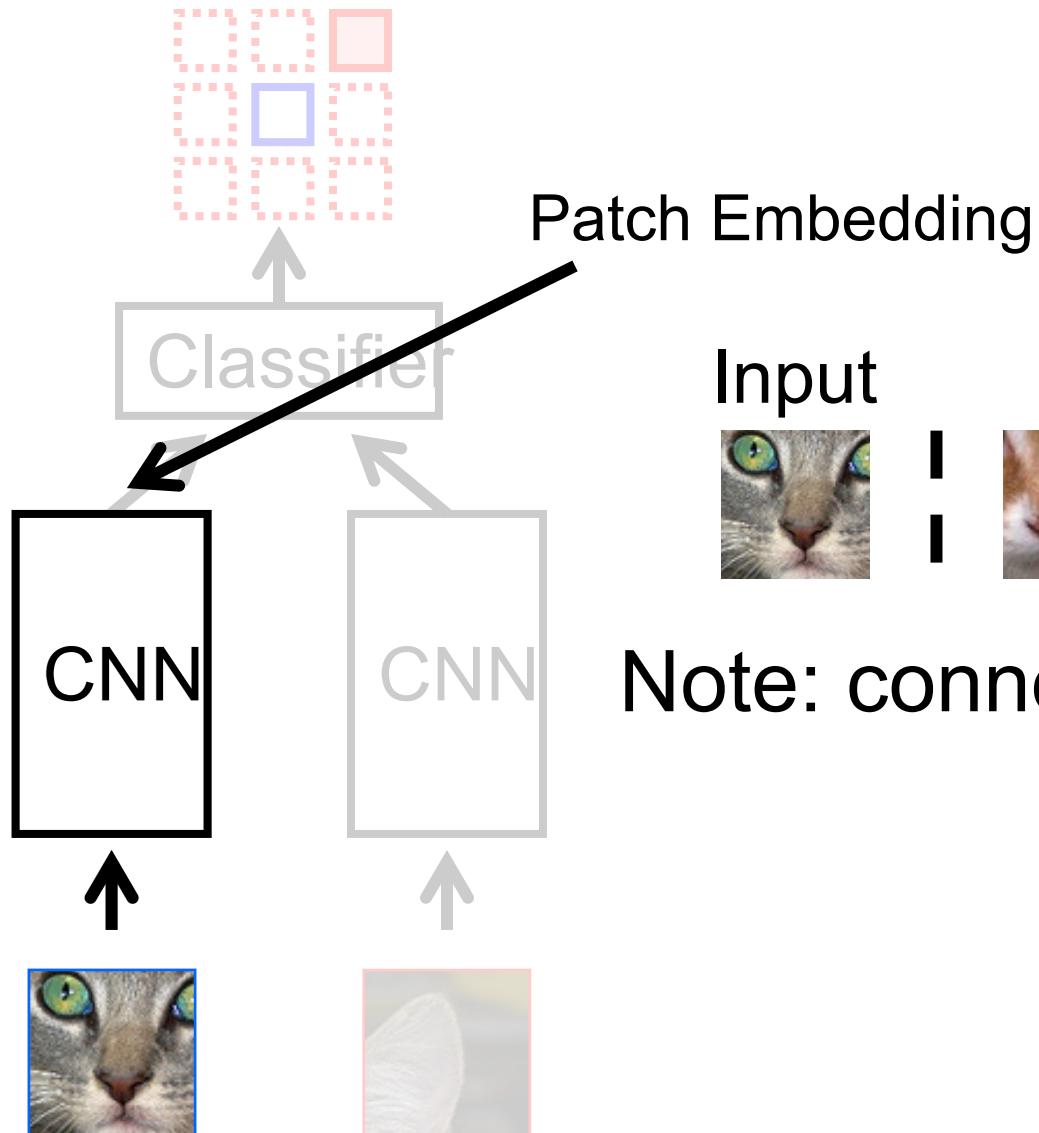
?

Semantics from a non-semantic task



Relative Position Task

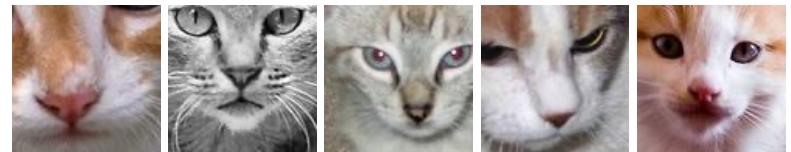




Input

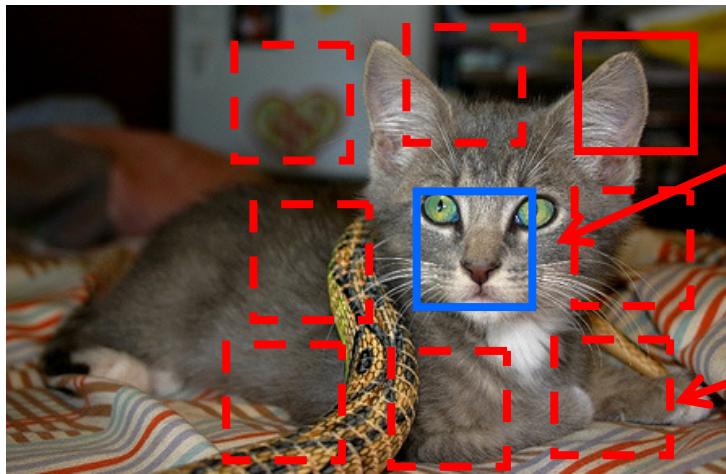
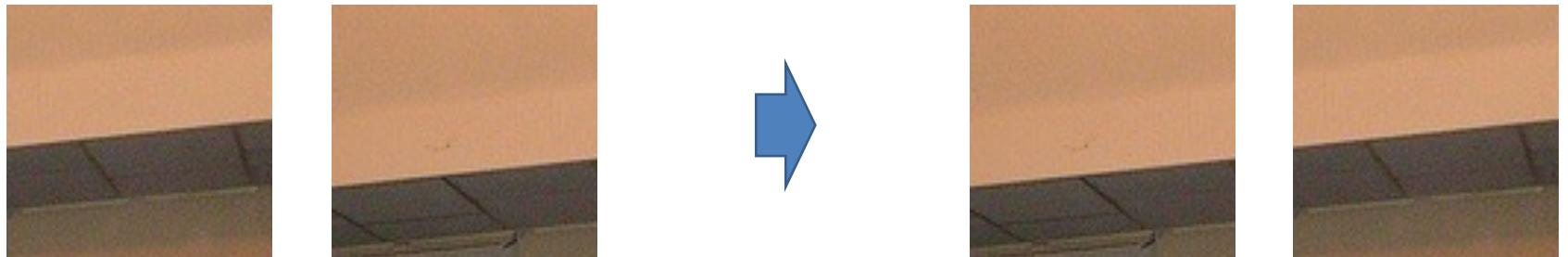


Nearest Neighbors



Note: connects **across** instances

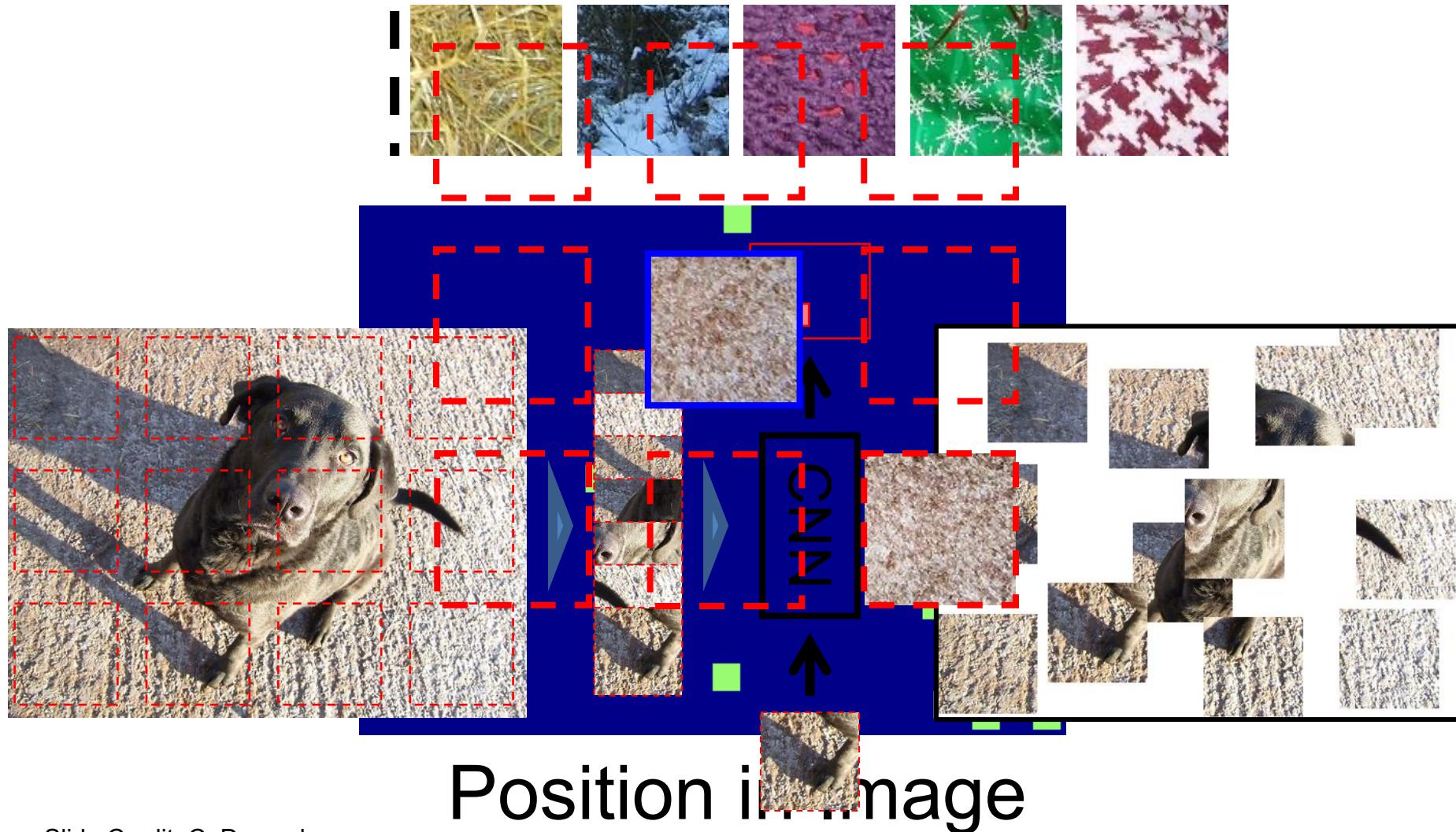
Avoiding Trivial Shortcuts



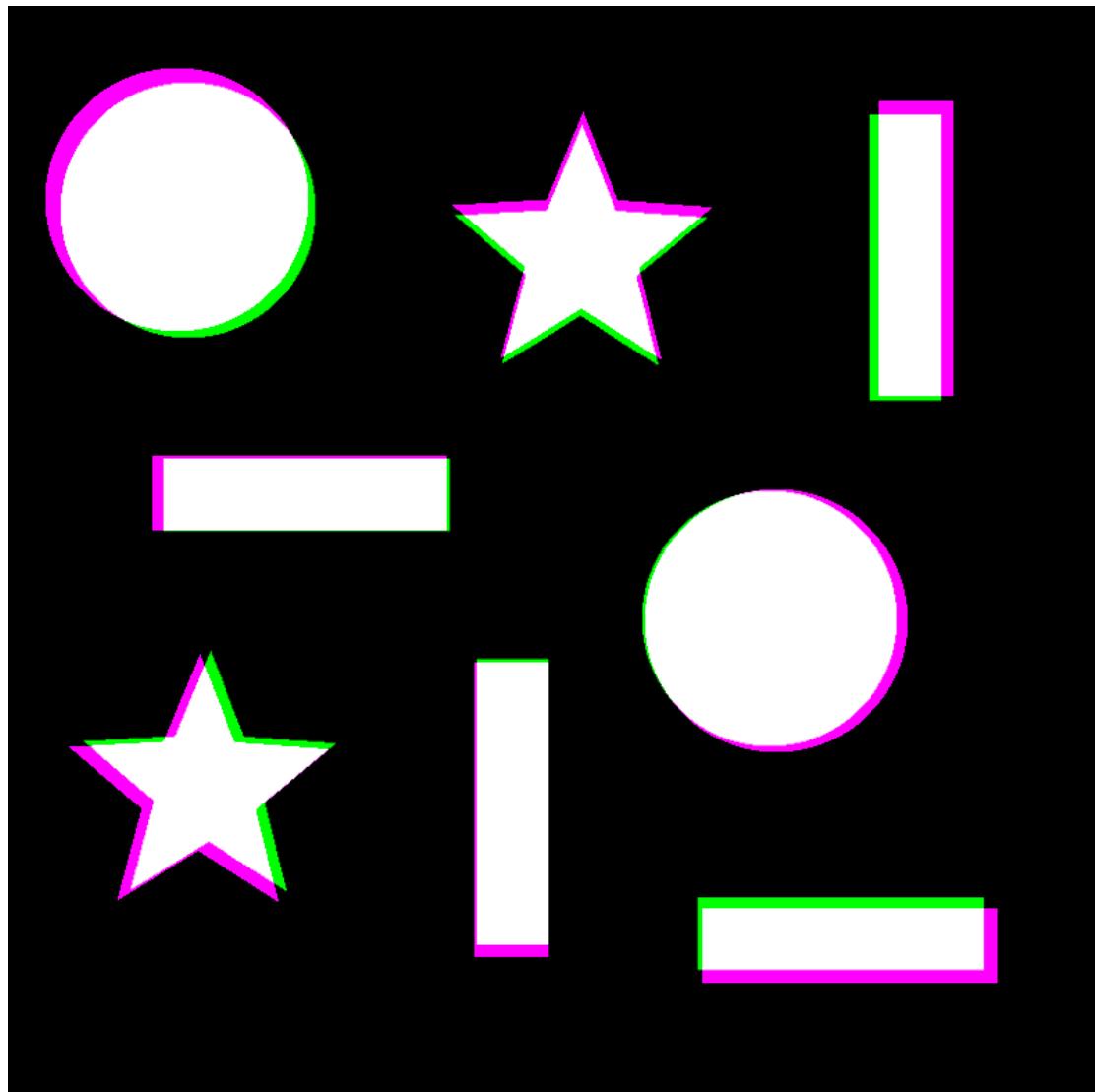
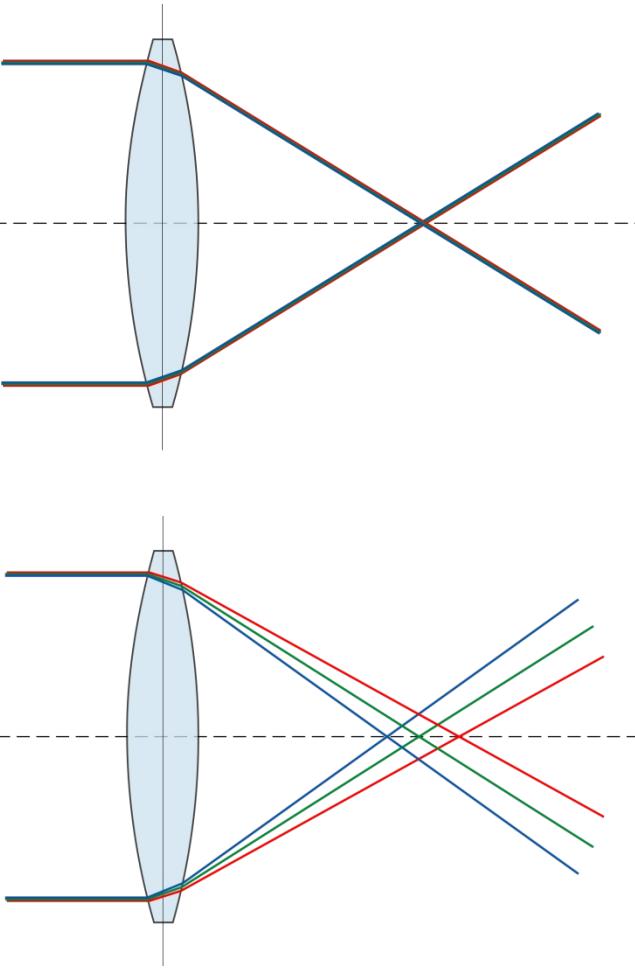
Include a gap

Jitter the patch locations

A Not-So “Trivial” Shortcut

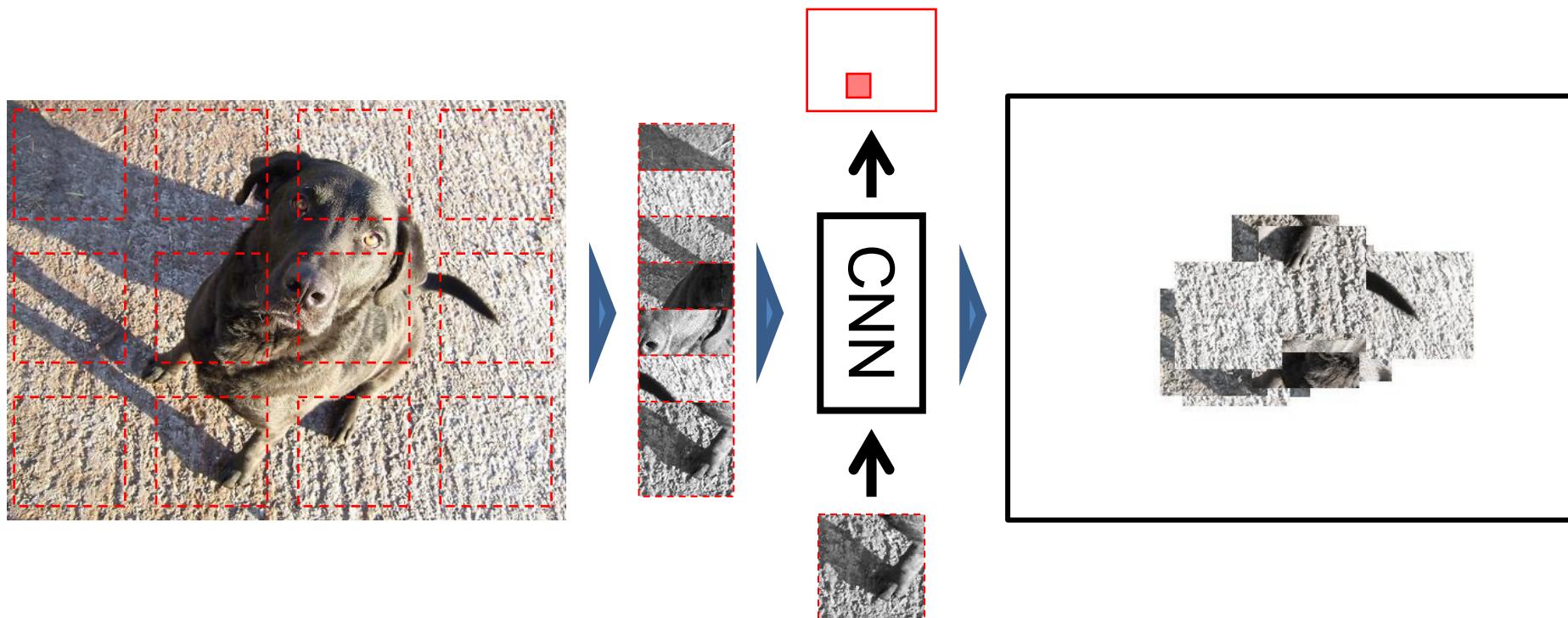


Chromatic Aberration

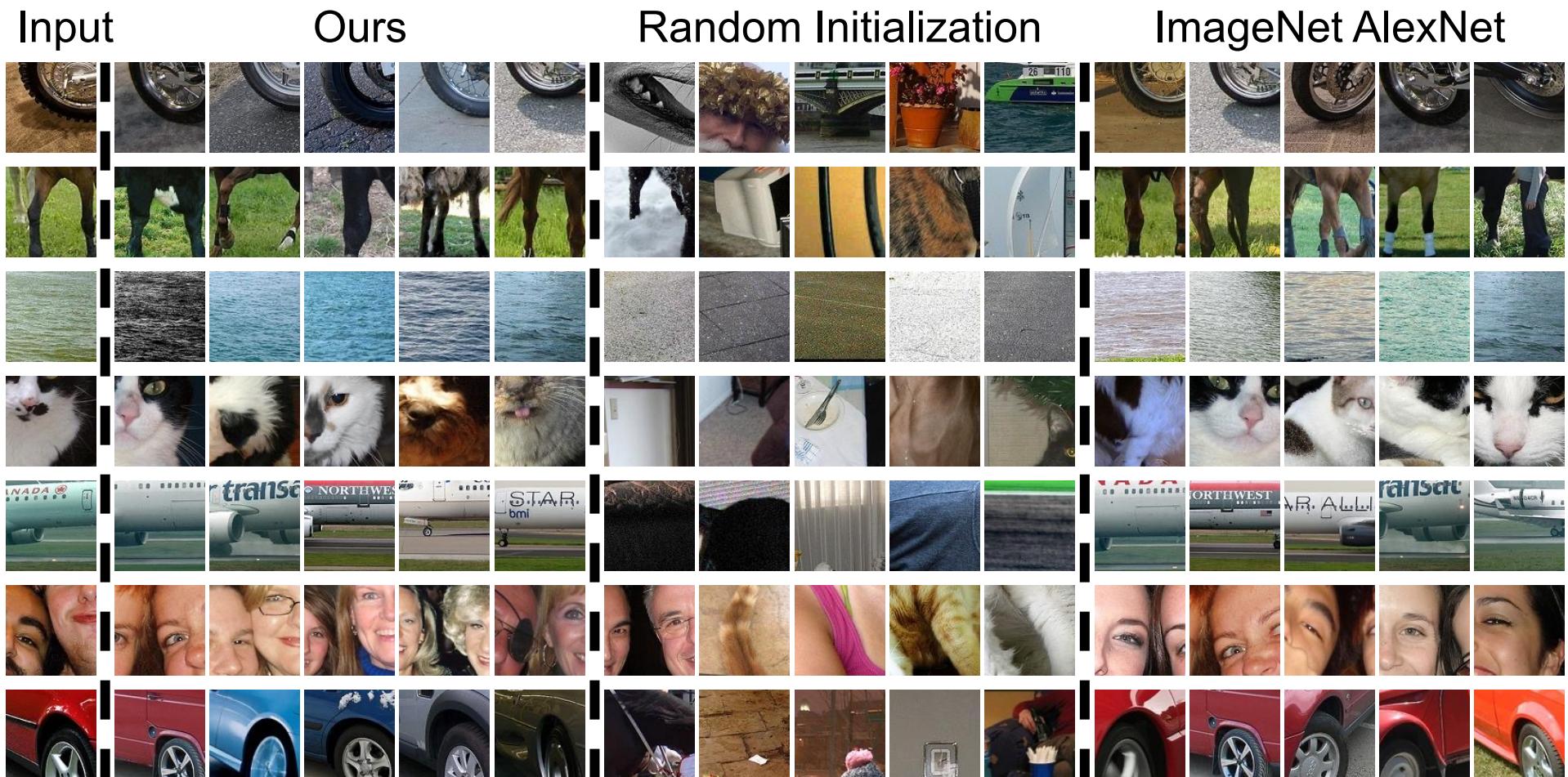


Slide Credit: C. Doersch

Chromatic Aberration



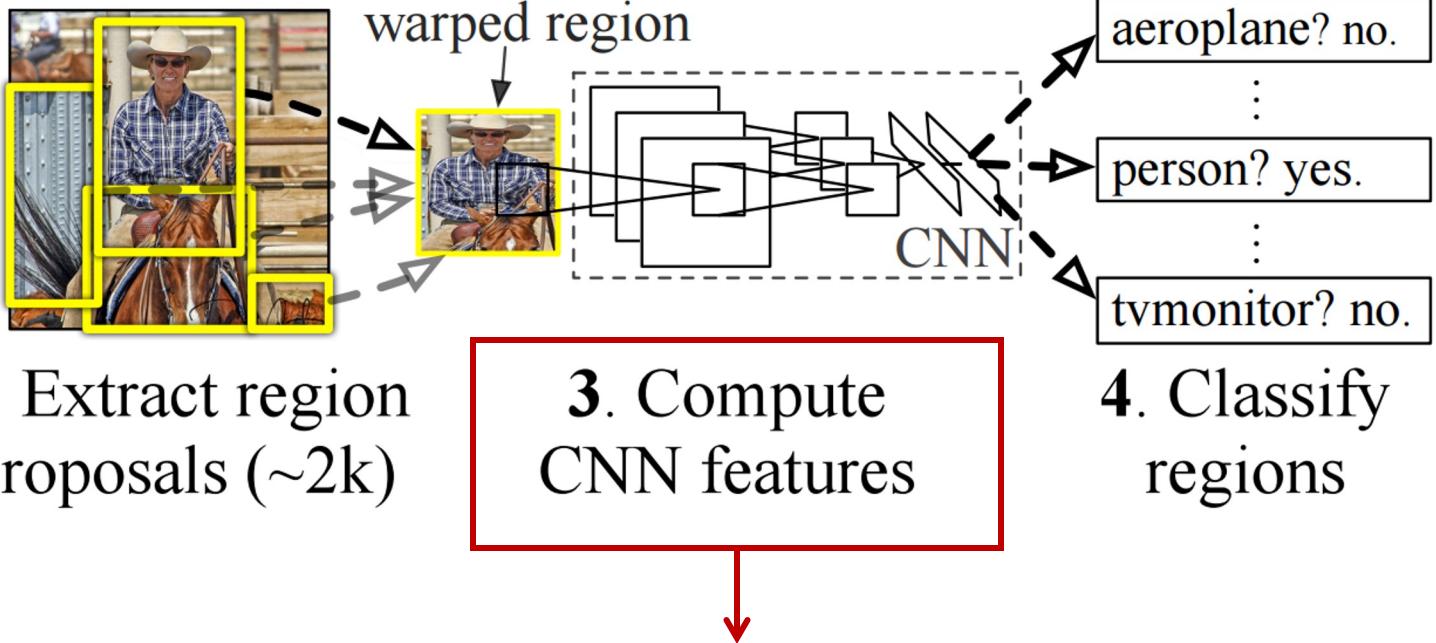
What is learned?



Pre-Training for R-CNN

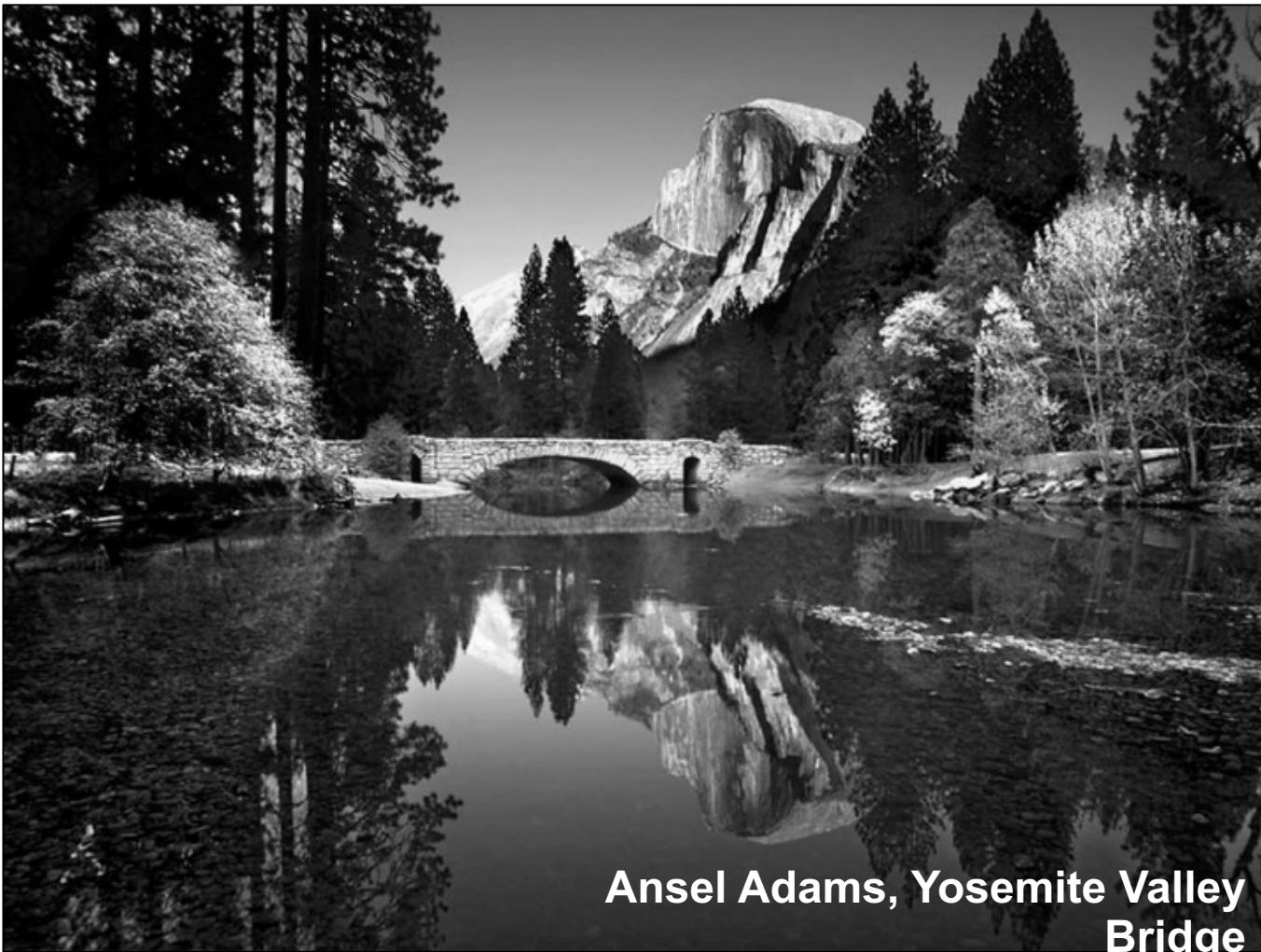


1. Input image



Pre-train on relative-position task, w/o labels

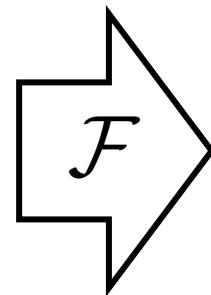
Other Sources Of Signal



**Ansel Adams, Yosemite Valley
Bridge**



Ansel Adams, Yosemite Valley Bridge – Our Result

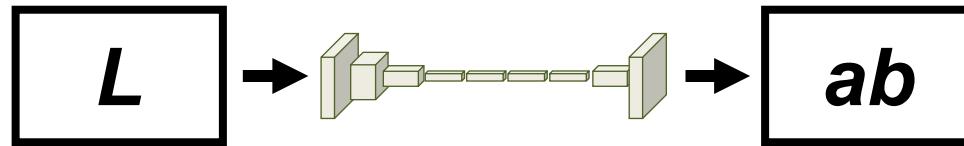


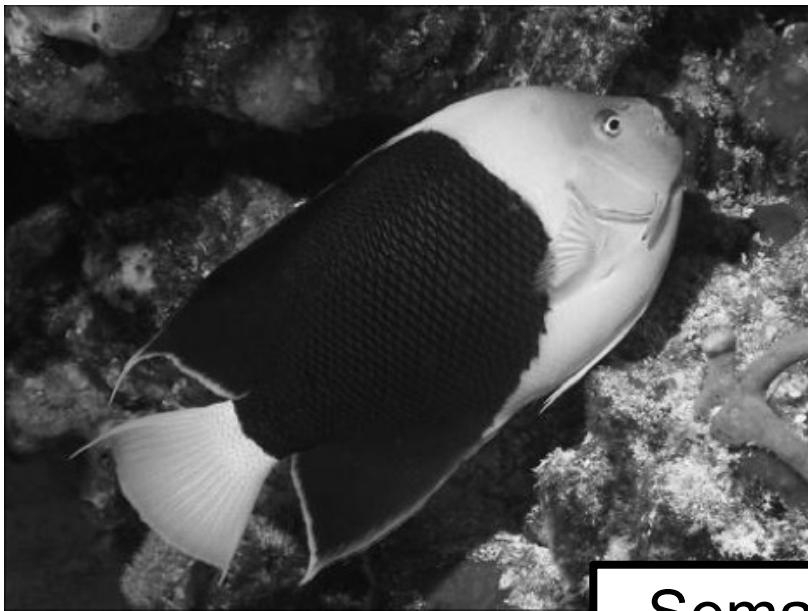
Grayscale image: L channel

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Color information: ab channels

$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$





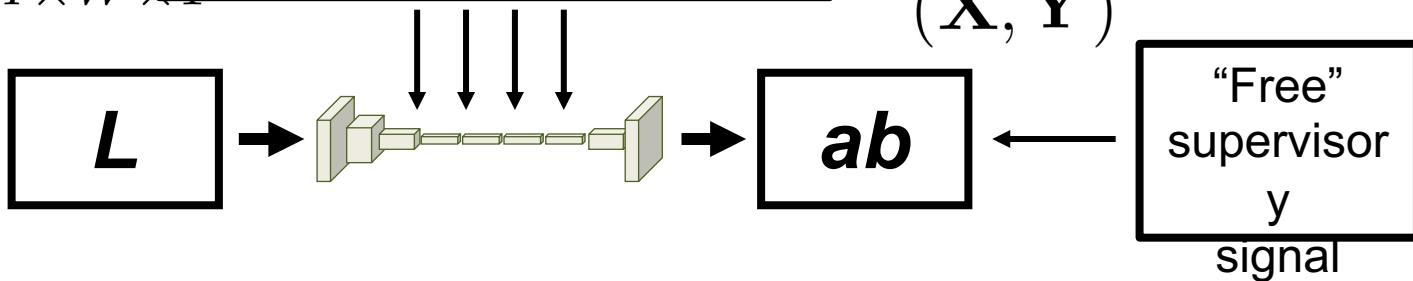
$$\xrightarrow{\mathcal{F}}$$



Grayscale image: L
 $X \in \mathbb{R}^{H \times W \times 1}$

Semantics? Higher-level abstraction?

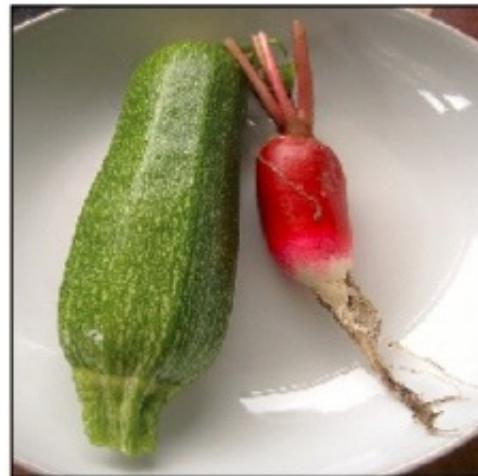
Concatenate (L, ab)
(X, \hat{Y})



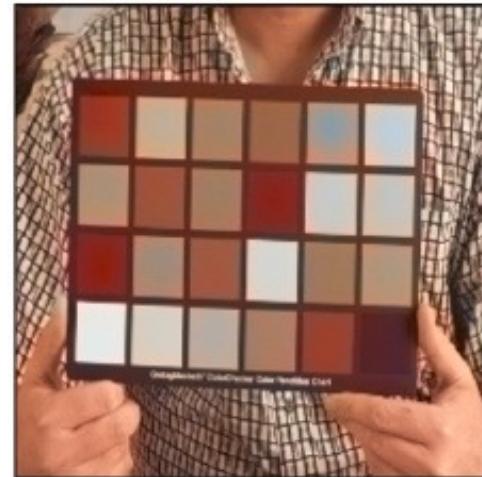
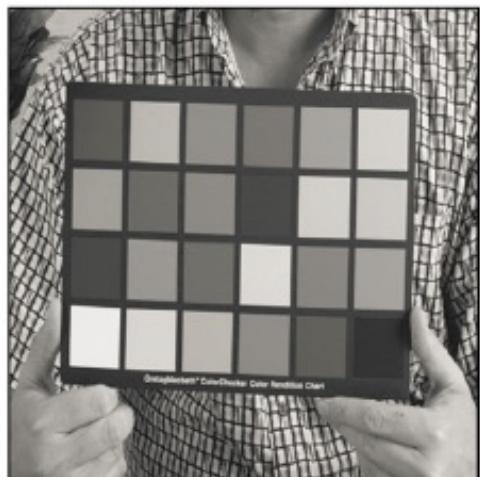
Input



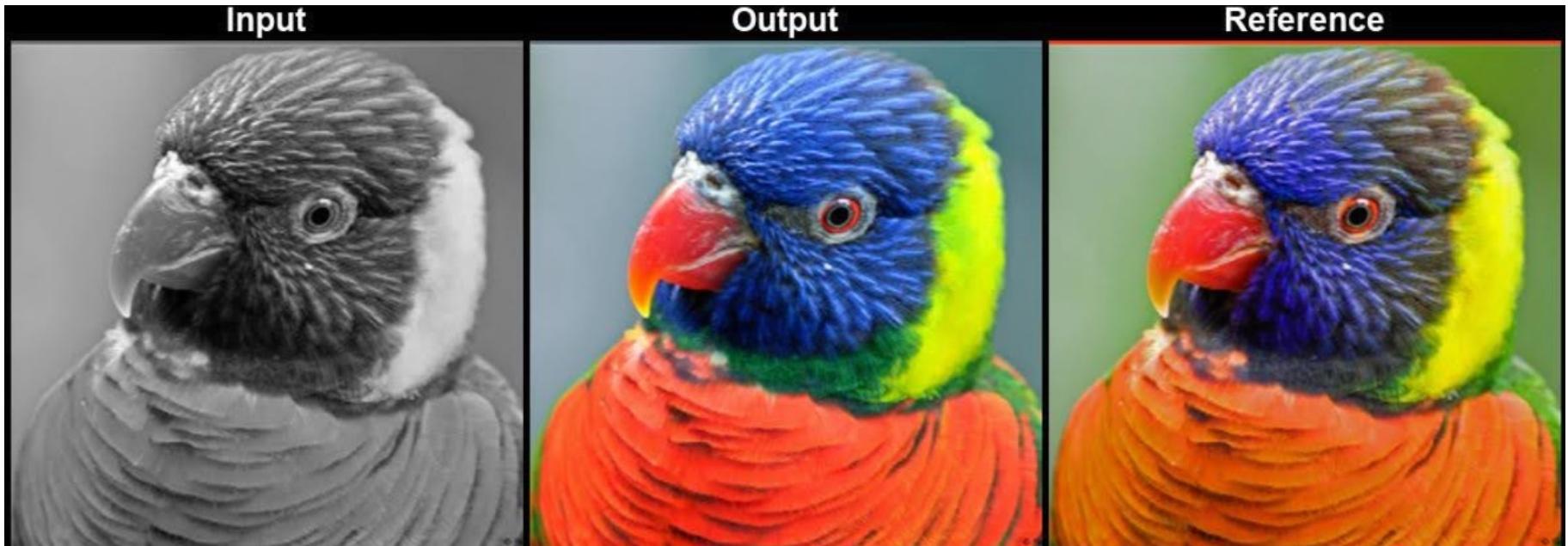
Ground Truth



Output



Got Much Better with Diffusion Models





Visually Indicated Sounds

Andrew Owens

Phillip Isola

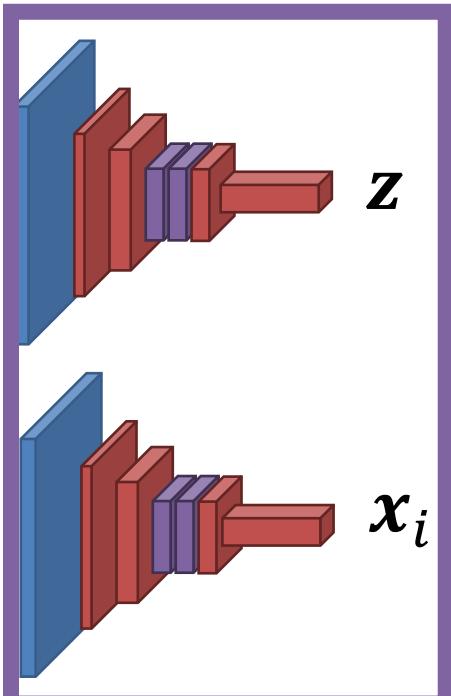
Josh McDermott

Antonio Torralba

Edward Adelson

William Freeman

Contrastive Learning



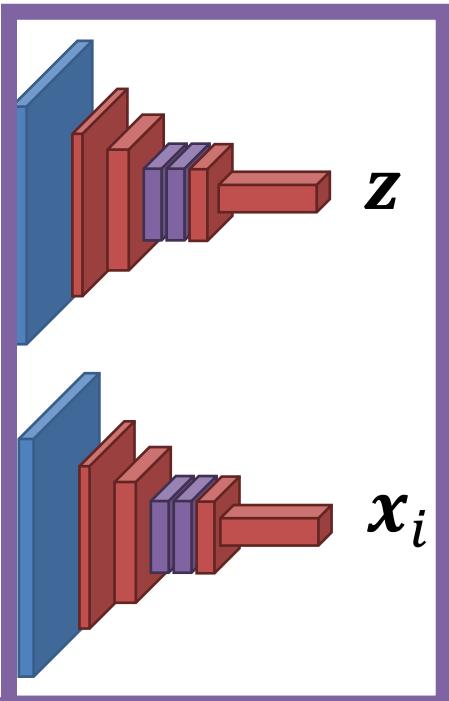
Weights shared
between networks

Given sample, construct feature z ; take a bunch of other images, minimize:

$$-\log \left(\frac{\exp(\mathbf{z}^T \mathbf{z})}{\exp(\mathbf{z}^T \mathbf{z}) + \sum_i \exp(\mathbf{z}^T \mathbf{x}_i)} \right)$$

$$\frac{\exp(\mathbf{z}^T \mathbf{z})}{\exp(\mathbf{z}^T \mathbf{z}) + \sum_i \exp(\mathbf{z}^T \mathbf{x}_i)} \leq 1$$

Contrastive Learning



Weights shared
between networks

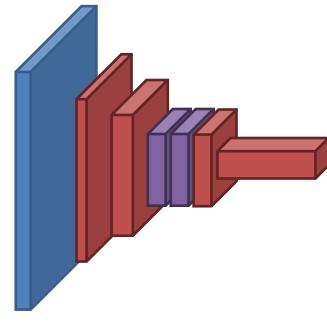
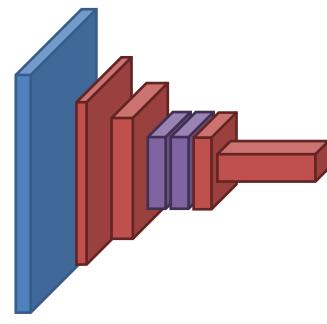
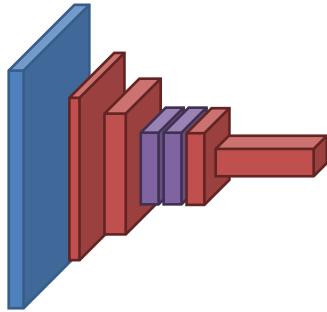
Given sample, construct feature z ; take a bunch of other images, minimize:

$$-\log \left(\frac{\exp(z^T z)}{\exp(z^T z) + \sum_i \exp(z^T x)} \right)$$

Basically, a scoring function with $w = z$ to measure similarity:

$$\frac{\exp(w^T z)}{\exp(w^T z) + \sum_i \exp(w^T x)}$$

Contrastive Learning



Best performing methods measure distance to augmented sample \mathbf{z}' :

$$-\log \left(\frac{\exp(\mathbf{z}^T \mathbf{z}')}{\exp(\mathbf{z}^T \mathbf{z}') + \sum_i \exp(\mathbf{z}^T \mathbf{x}_i)} \right)$$

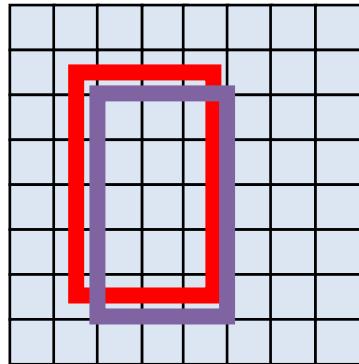
Goal: score augmented sample higher than everything else.

Next Time

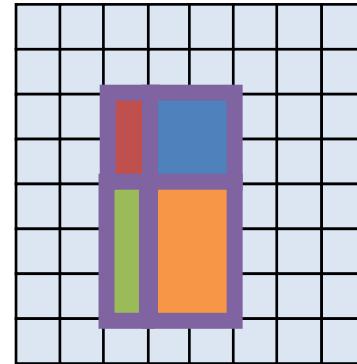
- Synthesizing Images

Extra Stuff

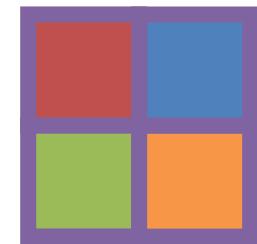
Fast R-CNN – ROI-Pool



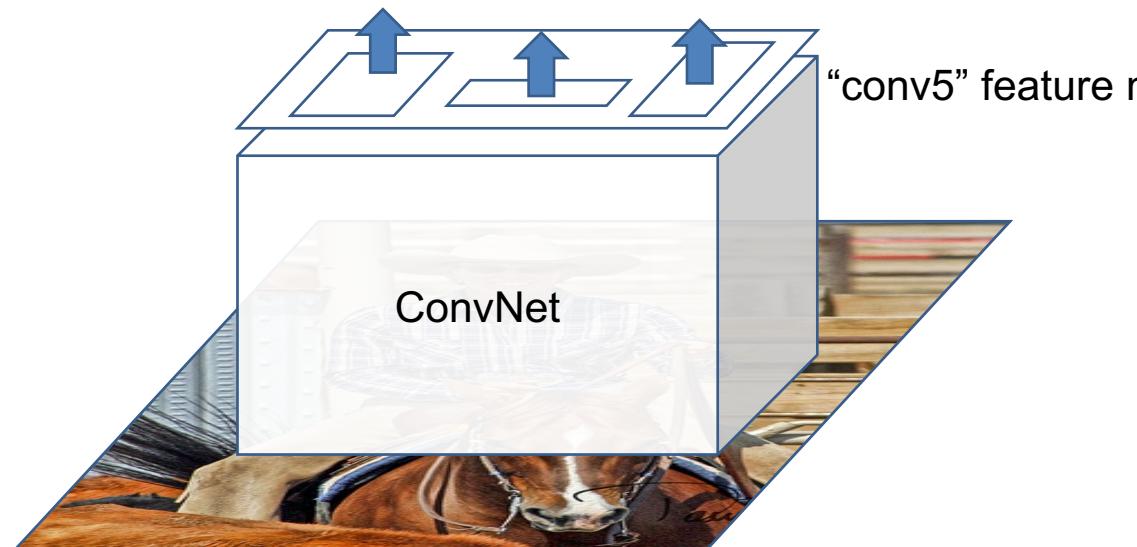
Line up



Divide

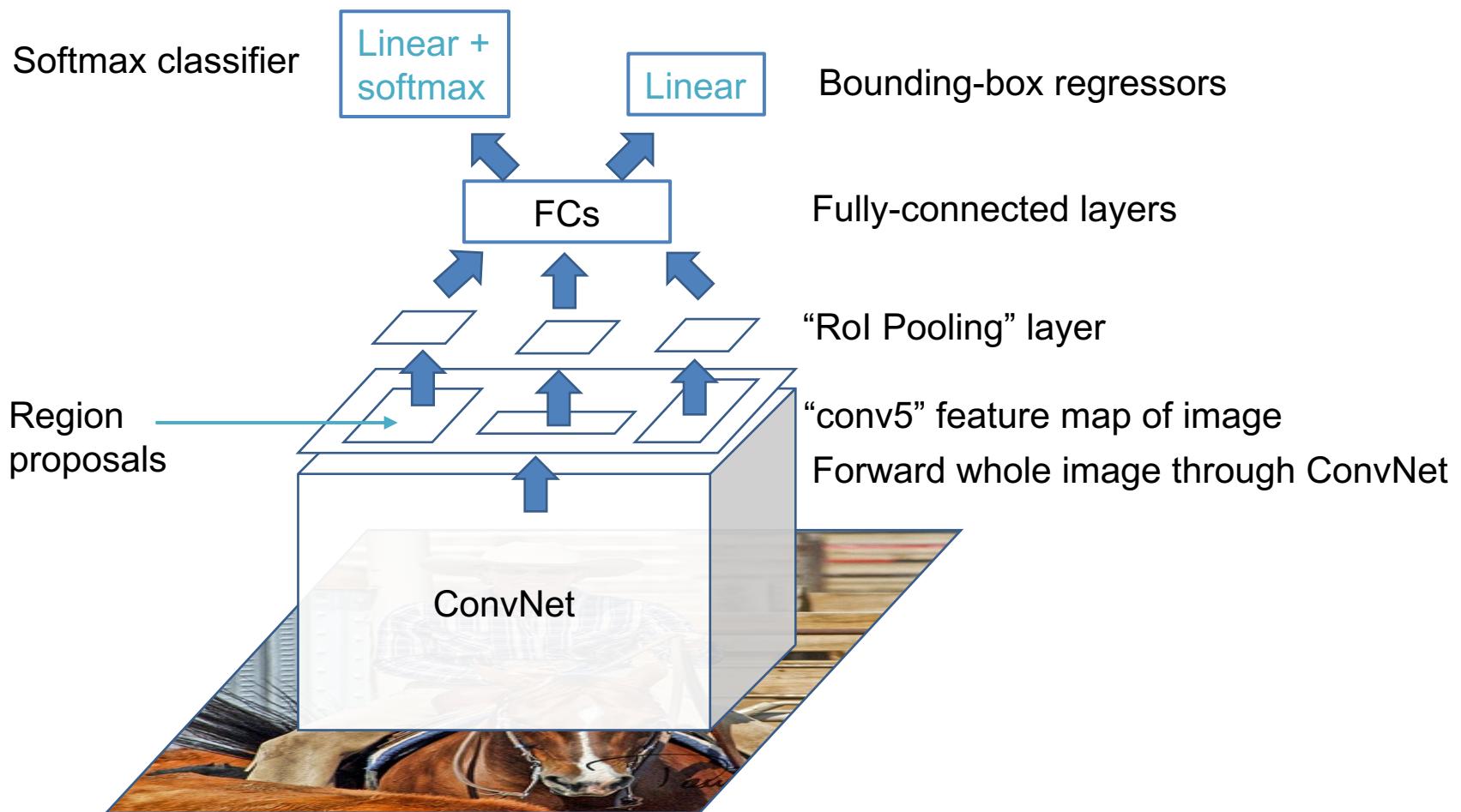


Pool

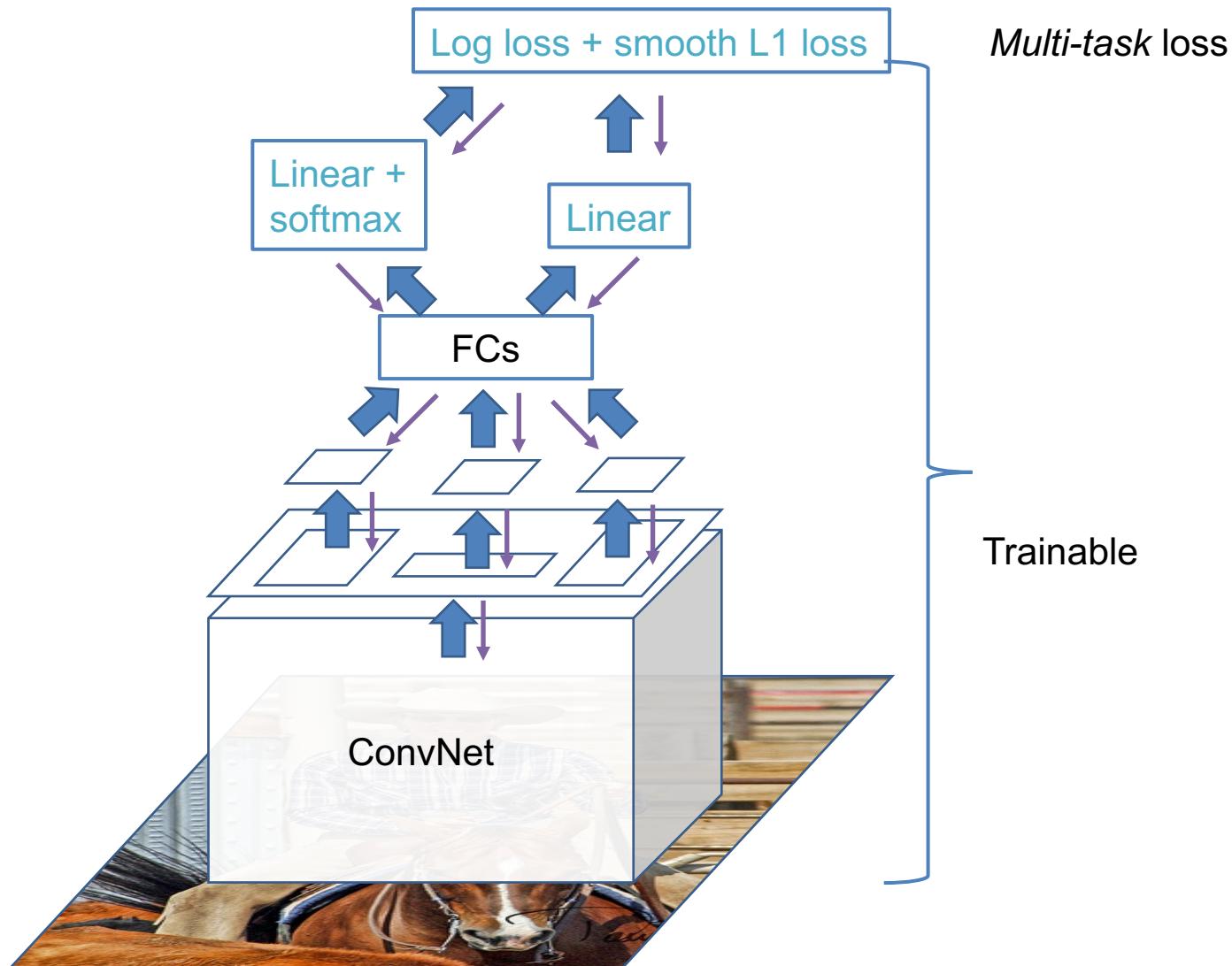


“conv5” feature map of image

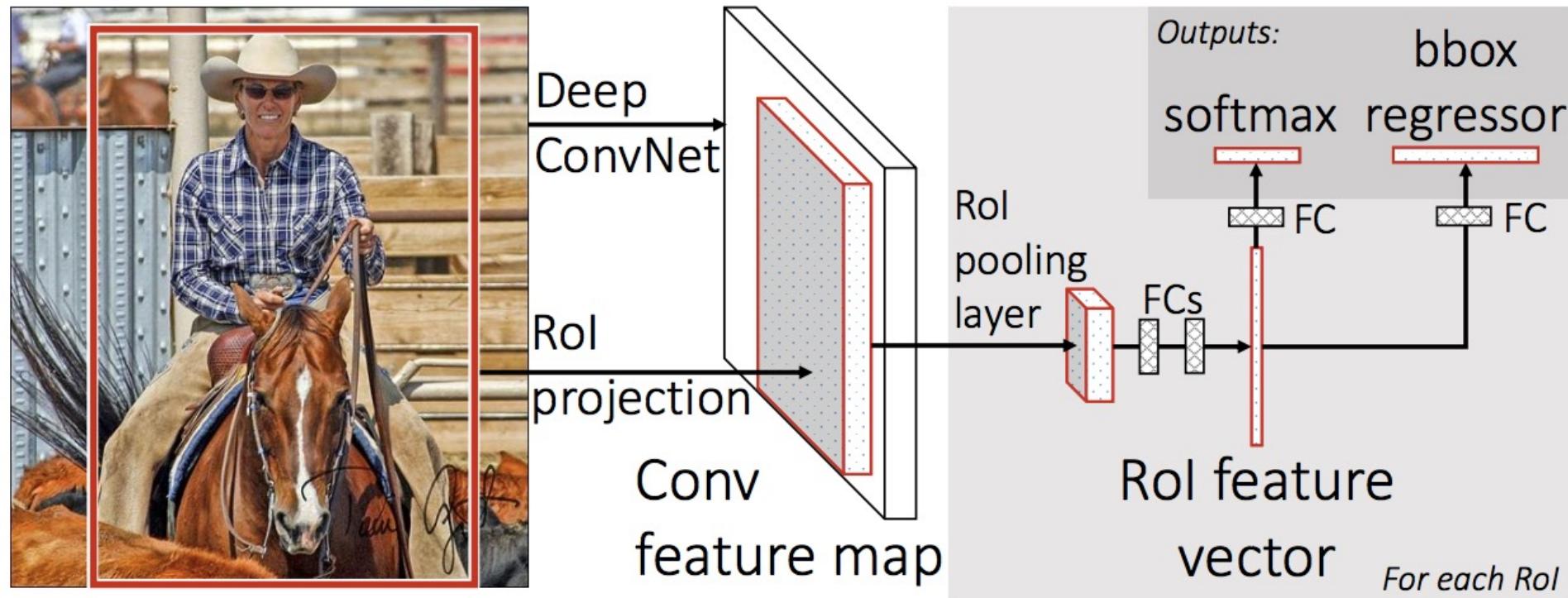
Fast R-CNN



Fast R-CNN training



Fast R-CNN: Another view



Fast R-CNN results

	Fast R-CNN	R-CNN
Train time (h)	9.5	84
- Speedup	8.8x	1x
Test time / image	0.32s	47.0s
Test speedup	146x	1x
mAP	66.9%	66.0%

Timings exclude object proposal time, which is equal for all methods.
All methods use VGG16 from Simonyan and Zisserman.